

# 7.91 – Lecture #1 Michael Yaffe

## Introduction to Bioinformatics

*Focus on Kinases*

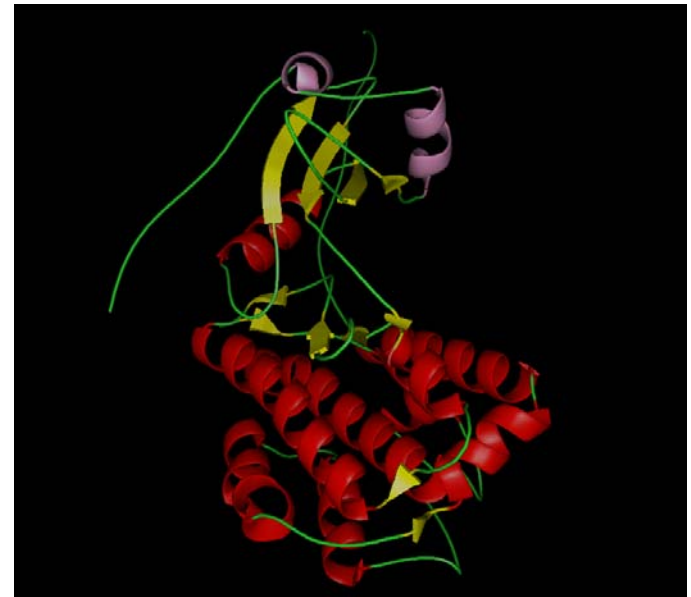
## & Pairwise Sequence Comparisons

```
ARDFSHGLENKLLGCDSMRWE
. . . . .
GRDYKMALLEQWILGCD-MRWD
```

Reading:

This lecture: Mount pp.1-7, 29-35, 45-48, 51-64

Next lecture: Mount pp. 8-9, 65-89, 96-115, +more



The Protein Data Bank (PDB - <http://www.pdb.org/>) is the single worldwide repository for the processing and distribution of 3-D biological macromolecular structure data.

Berman, H. M., J. Westbrook, Z. Feng, G. Gilliland, T. N. Bhat, H. Weissig, I. N. Shindyalov, and P. E. Bourne. The Protein Data Bank. *Nucleic Acids Research* 28 (2000): 235-242.

(PDB Advisory Notice on using materials available in the archive: <http://www.rcsb.org/pdb/advisory.html>)

# Outline

- Review of biological fundamentals
- Genetics example
- Biological databases/NCBI resources
- Simplified sequence analysis – where to start...
- Simple sequence comparisons
- Definitions of related sequences
- Concepts and types of alignments – the good, the bad, and the ugly
- Dot matrix alignments
- Computational efficiency
- Recursion and dynamic programming
- Substitution matrices: PAM, BLOSUM, Gonnet

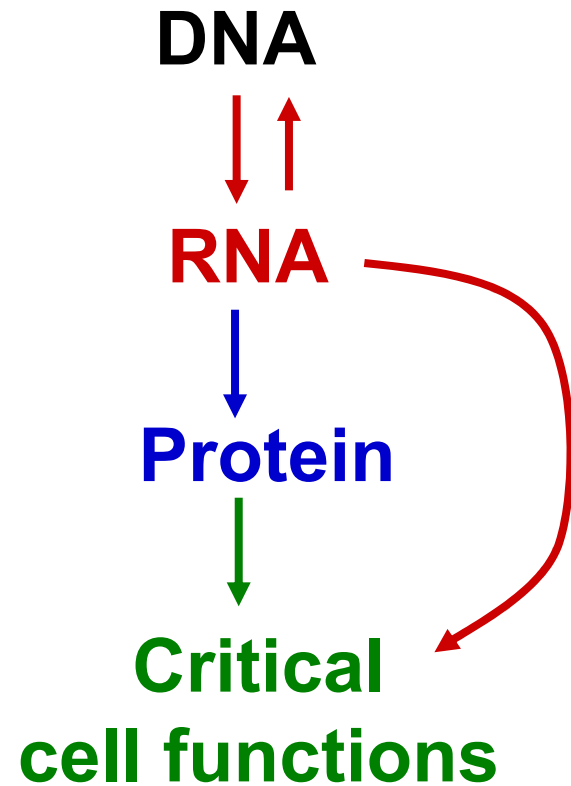
# Outline (cont)

- Gaps
- Applied dynamic programming: global alignments: Needleman-Wunsch
- Applied dynamic programming: local alignments – Smith-Waterman
- Basic statistics of sequence alignments

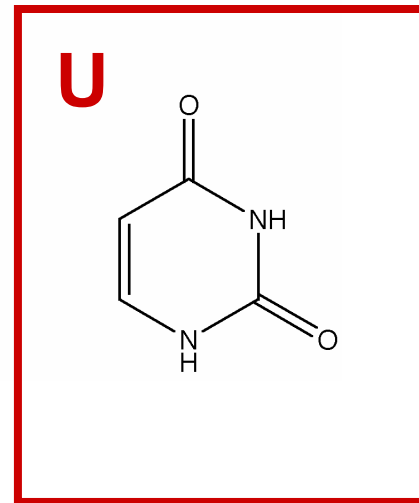
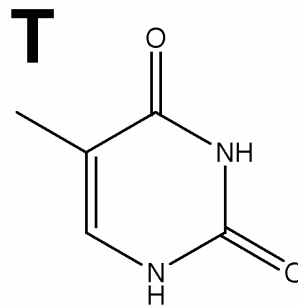
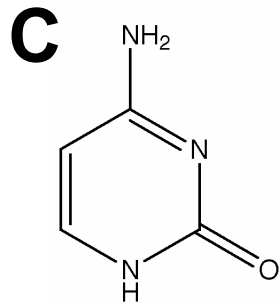
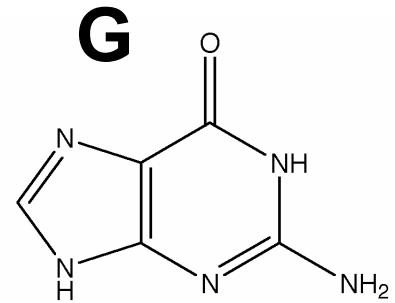
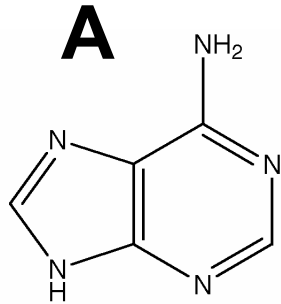
# Review of biological fundamentals

## Genetic material

- gene as a concept
- DNA as hereditary material

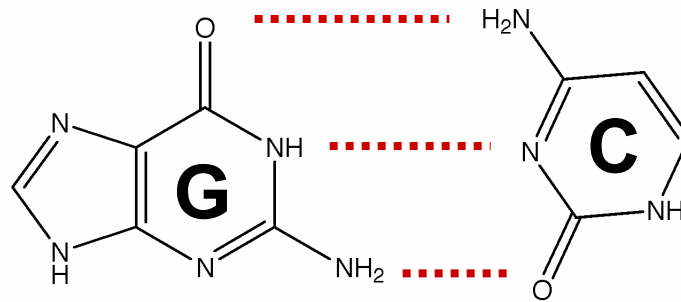
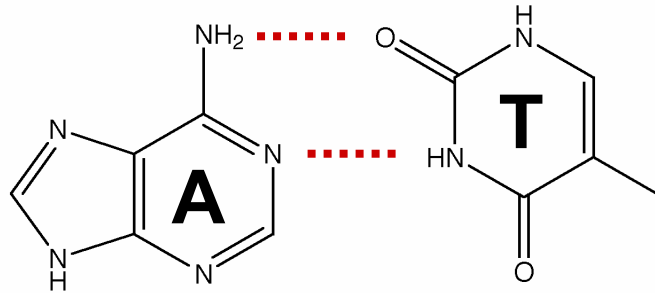


# DNA structure



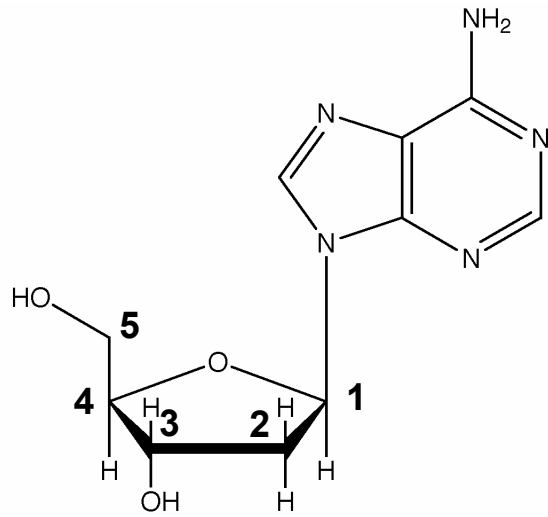
Bases

# DNA structure

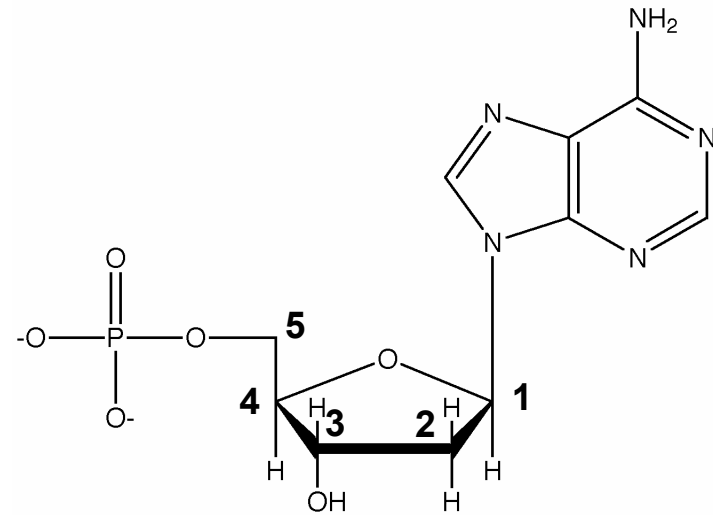


## Base pairing

# DNA structure

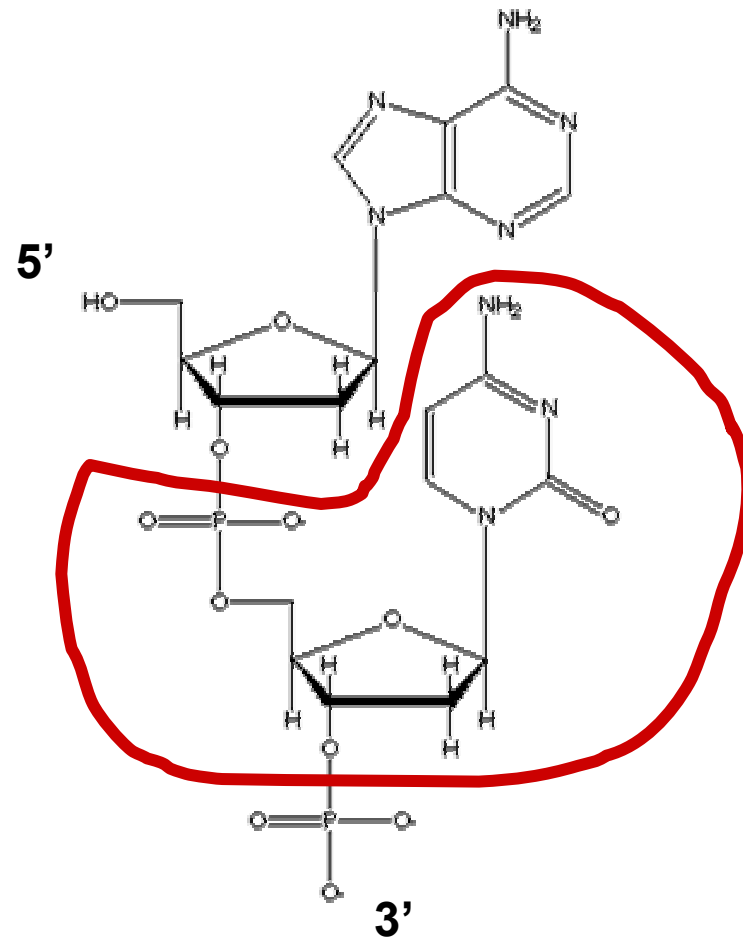


**Nucleoside**



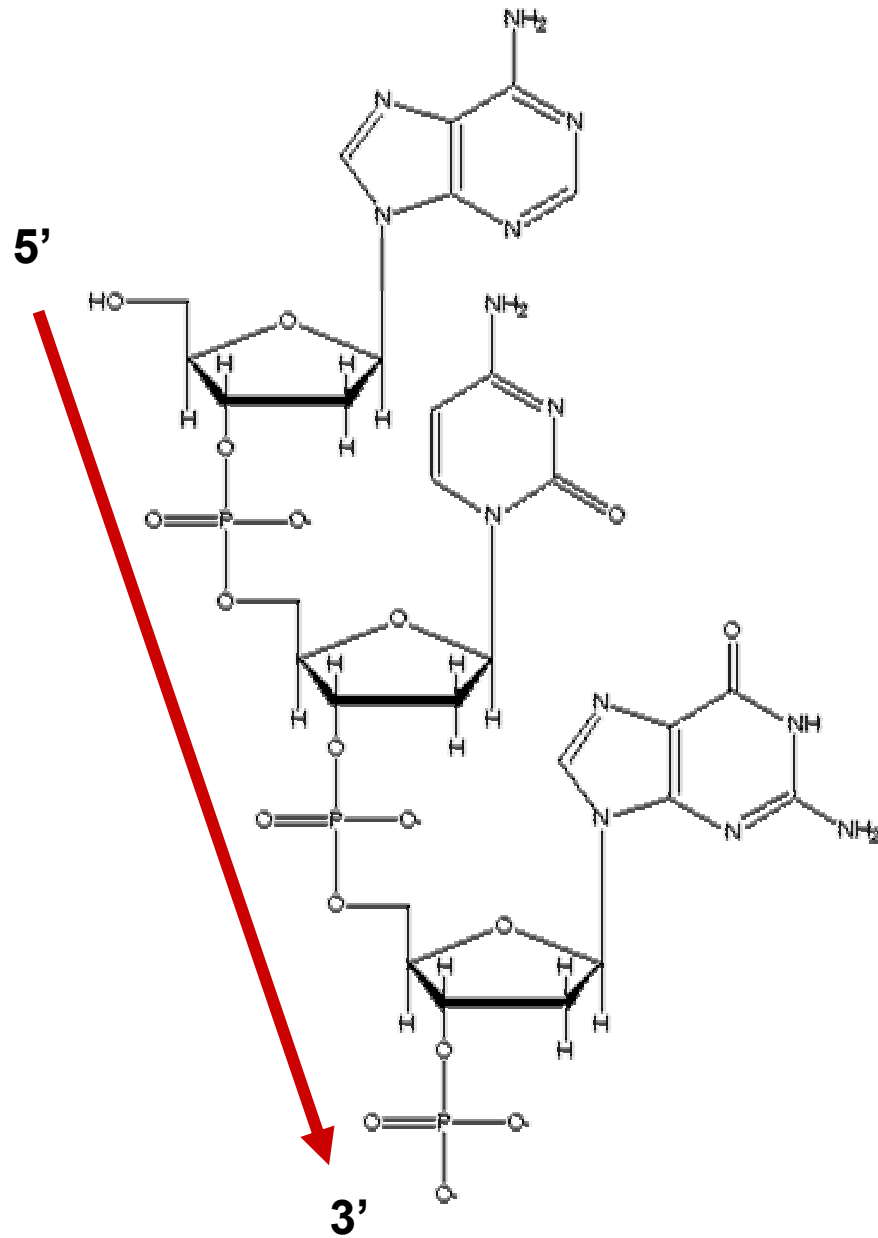
**Nucleotide**

# DNA structure

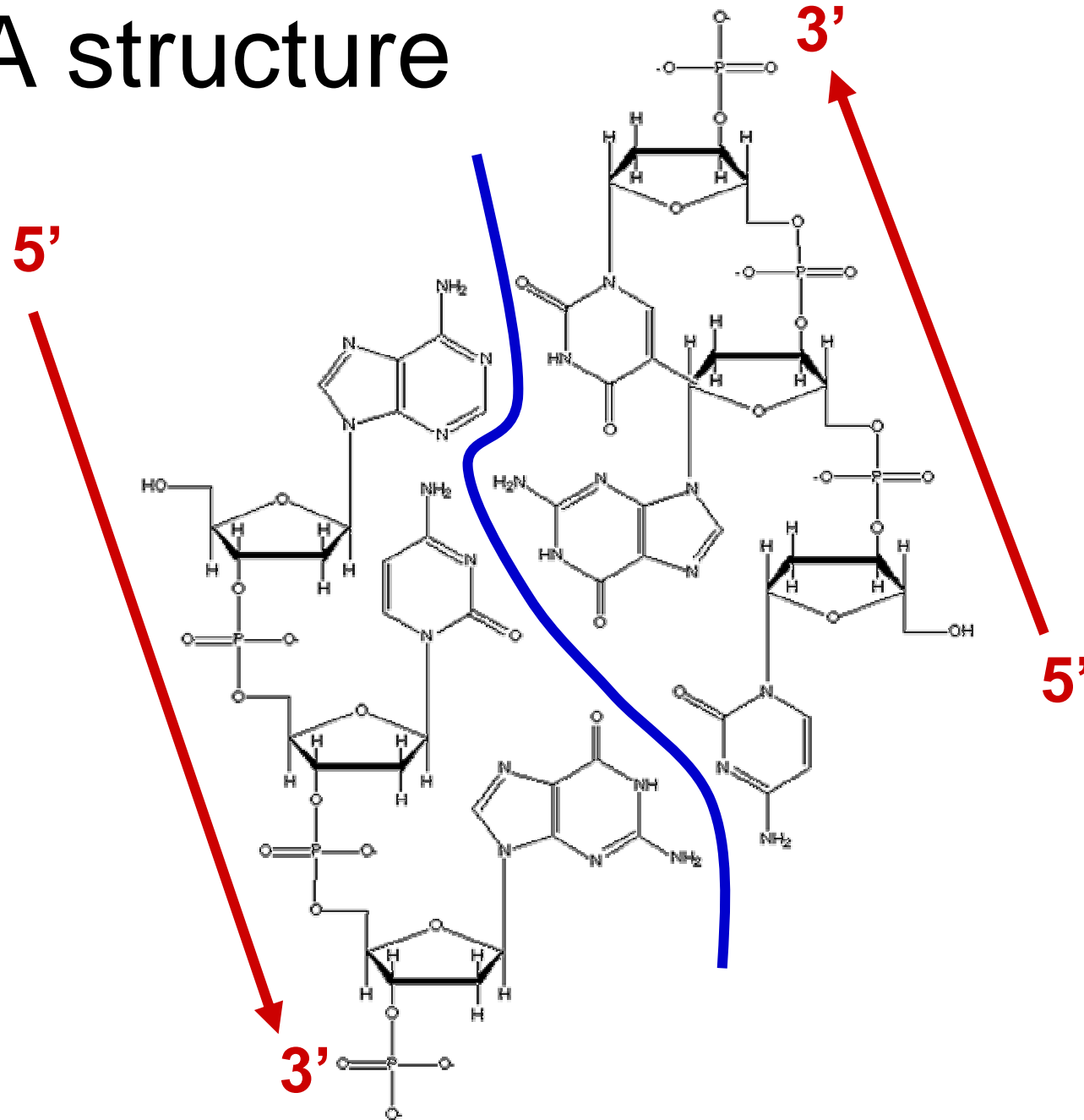




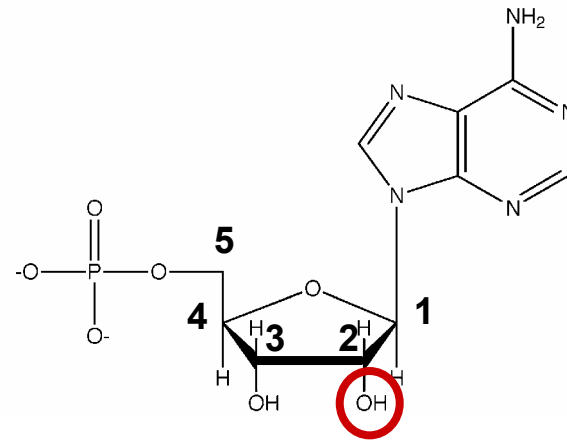
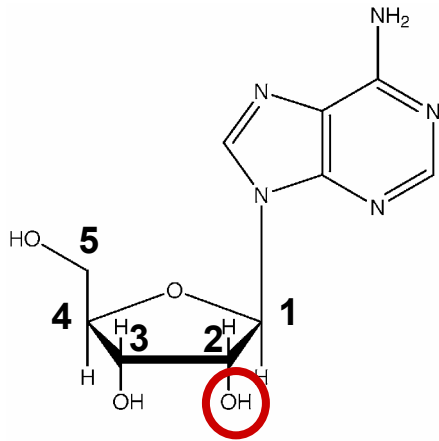
# DNA structure



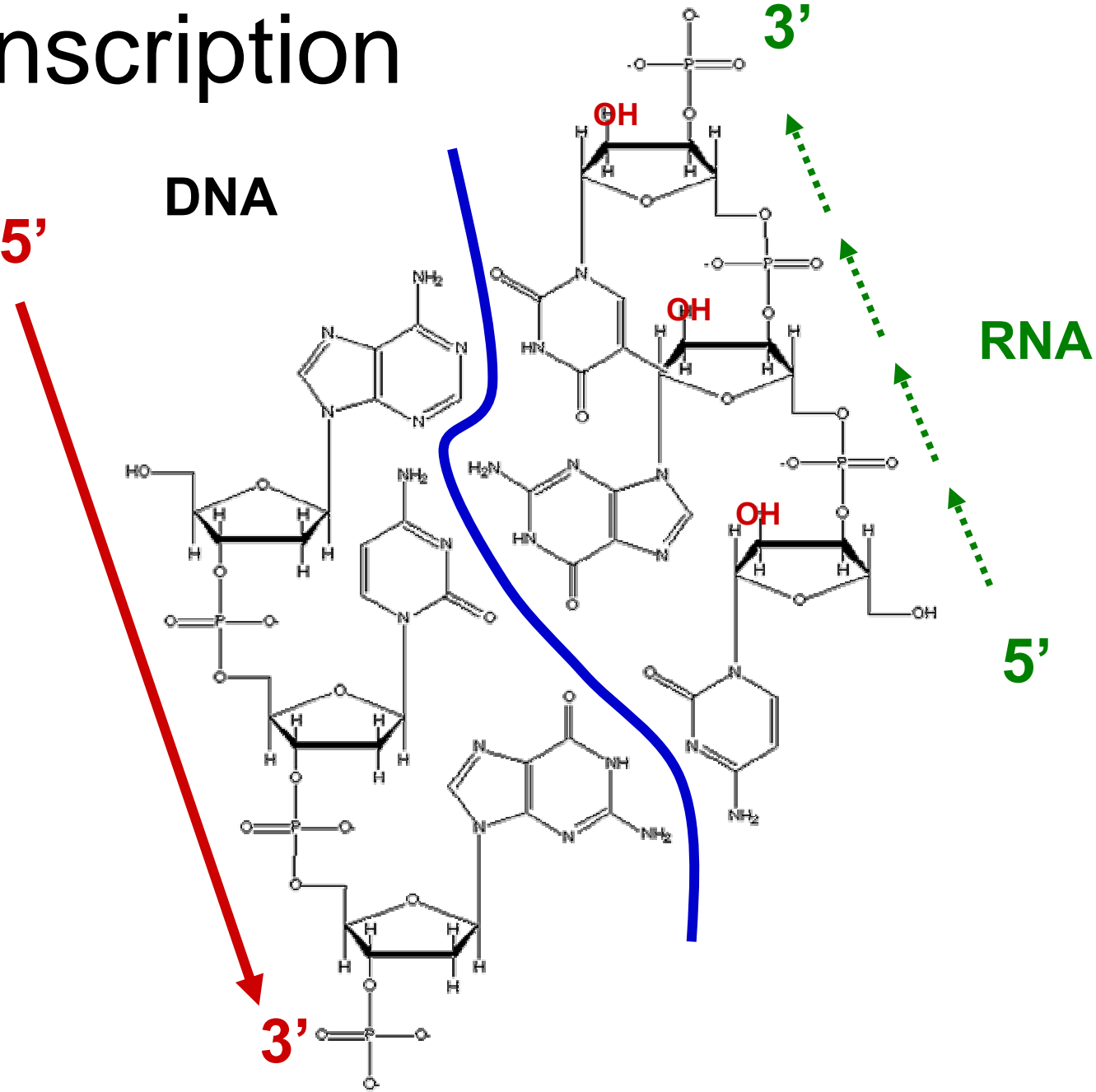
# DNA structure



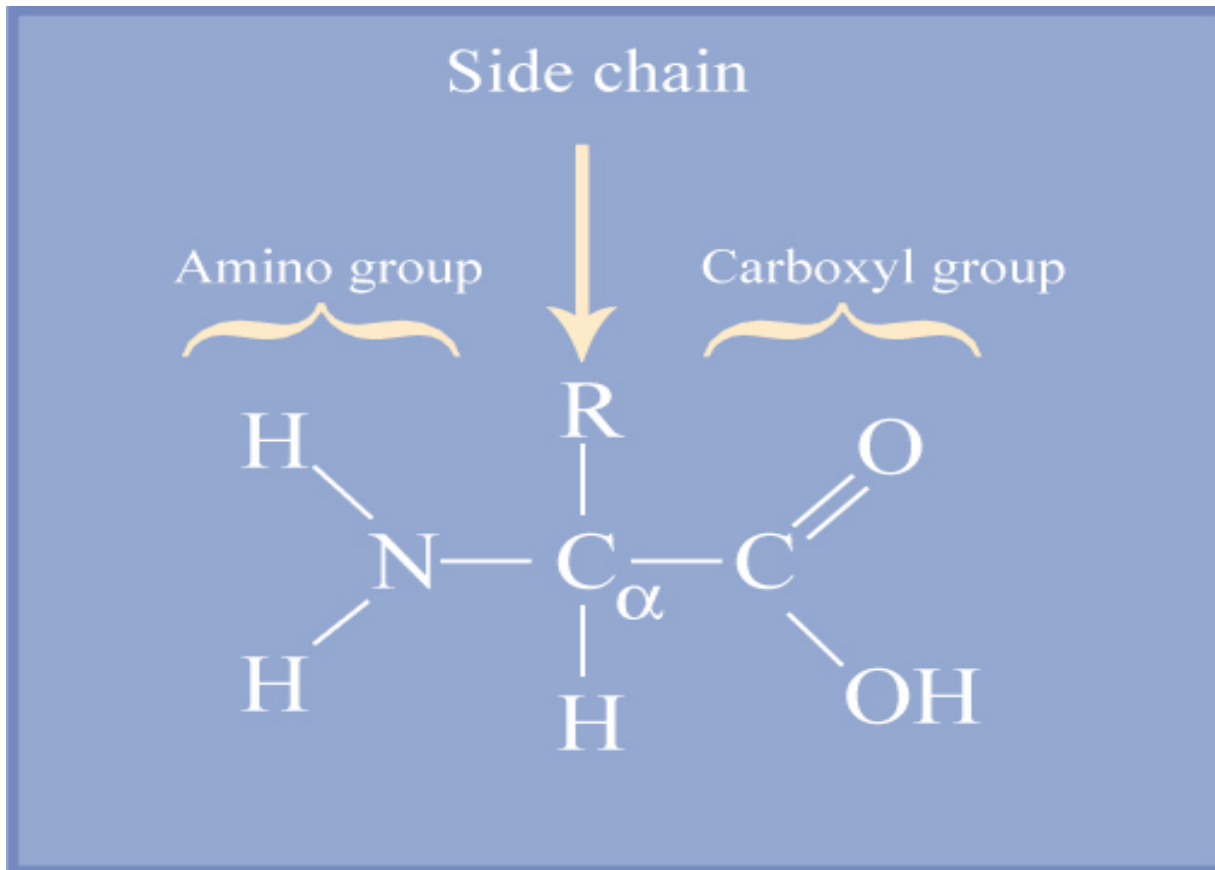
# RNA structure



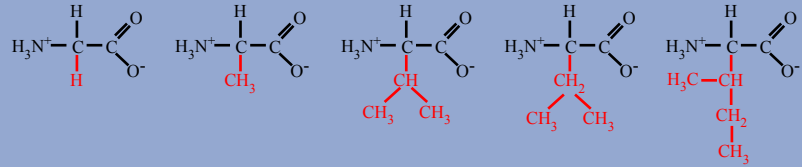
# Transcription



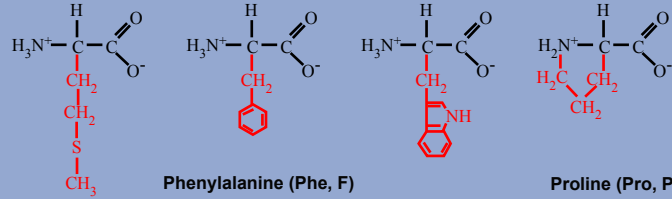
# Protein structure



Nonpolar, Hydrophobic R-groups

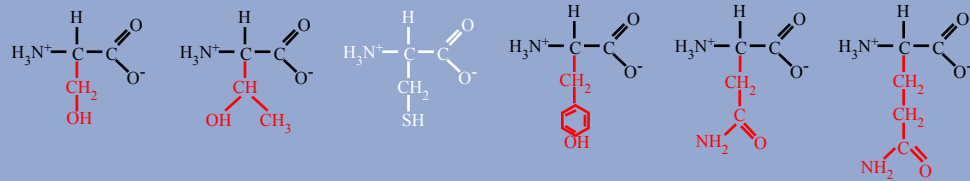


Glycine (Gly, G)    Alanine (Ala, A)    Valine (Val, V)    Leucine (Leu, L)    Isoleucine (Ile, I)



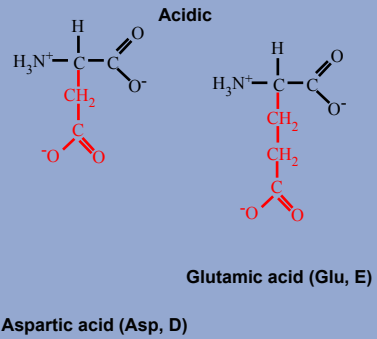
Methionine (Met, M)    Phenylalanine (Phe, F)    Tryptophan (Trp, W)    Proline (Pro, P)

Polar, Hydrophilic R-groups

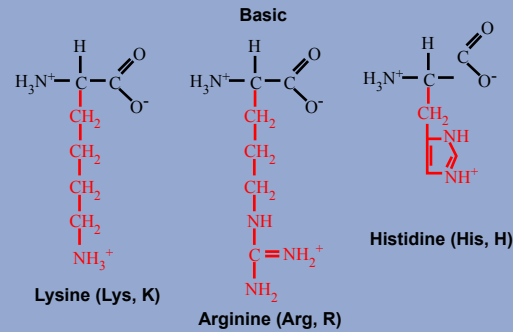


Serine (Ser, S)    Threonine (Thr, T)    Cysteine (Cys, C)    Tyrosine (Tyr, Y)    Asparagine (Asn, N)    Glutamine (Gln, Q)

Electrically charged



Aspartic acid (Asp, D)    Glutamic acid (Glu, E)



Lysine (Lys, K)    Arginine (Arg, R)    Histidine (His, H)

# Codon Table

(5')...pNpNpN...(3') in mRNA

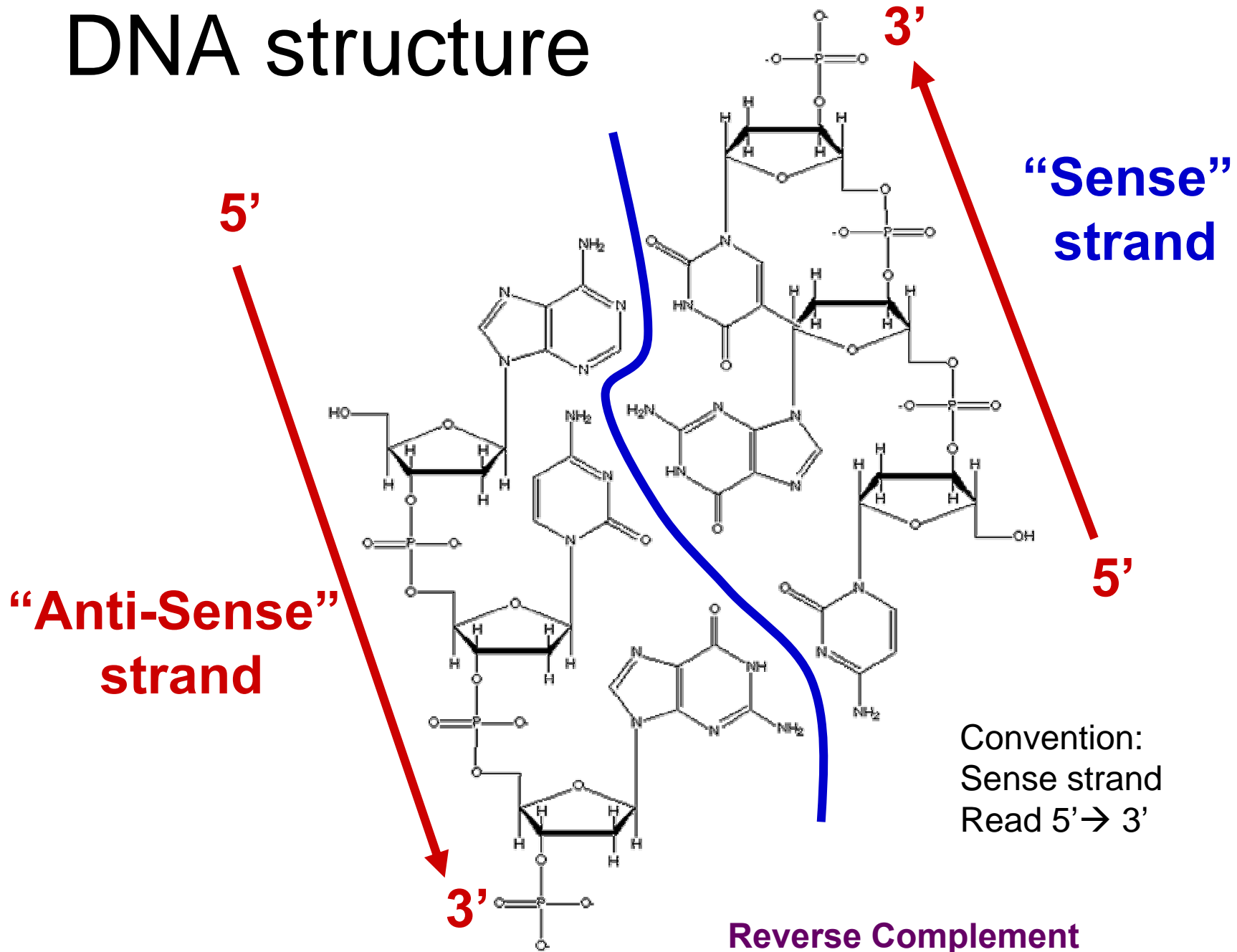
Middle Base of Codon →

Base at 5' End  
of Codon ↓

Base at 3' End  
of Codon ↓

|          | U                       | C   | A           | G           |   |
|----------|-------------------------|-----|-------------|-------------|---|
| <i>U</i> | phe (UUU)               | ser | tyr         | cys         | U |
|          | phe                     | ser | tyr         | cys         | C |
|          | leu                     | ser | termination | termination | A |
|          | leu                     | ser | termination | trp         | G |
| <i>C</i> | leu                     | pro | his         | arg         | U |
|          | leu                     | pro | his         | arg         | C |
|          | leu                     | pro | gln         | arg         | A |
|          | leu                     | pro | gln         | arg         | G |
| <i>A</i> | ile                     | thr | asn         | ser         | U |
|          | ile                     | thr | asn         | ser         | C |
|          | ile                     | thr | lys         | arg         | A |
|          | met<br>(and initiation) | thr | lys         | arg         | G |
| <i>G</i> | val                     | ala | asp         | gly         | U |
|          | val                     | ala | asp         | gly         | C |
|          | val                     | ala | glu         | gly         | A |
|          | val                     | ala | glu         | gly         | G |

# DNA structure





# Genetics experiment:

Isolate a yeast mutant that has increased chromosome number

phenotype

Can rescue the phenotype with a piece of DNA

...corresponds to a site of mutation in the yeasts DNA

genotype

```
CGTTTTCTGTAAAGCGCTTAATTTGTTTACCATTCTATAAAAACCTTGAGCTAAGGCCAACTGATGCA
ATTGCTCAAGTGAATGCATAAACAAAGCAAGATCATTCTTAGCGCAAAAAAACTGGGATTTTGAAATAC
AACAAAAGAAAGAAGTAAAAAGGGAATGCAACGCAATAGTTTAGTAAATATCAAACCTAAACGCTAATTCG
CCATCGAAAAAGACCACAACAAGACCAAATACGTCCAGGATCAATAAACCATGGAGAATATCCCATTCCG
CGCAGCAAAGAAACCCGAATTCAAAAATACCTTCACCTGTAAGAGAAAAATTGAACAGATTACCTGTAAA
CAATAAGAAGTTTTTTGGATATGGAAAGCTCCAAAATTCCATCACCTATAAGGAAAGCGACTTCTTCCAAA
ATGATACACGAAAATAAGAAGCTACCTAAATTTAAATCCCTATCACTCGATGACTTTGAACTGGGGAAGA
AATTAGGAAAGGGTAAATTCGGTAAAGTTTATTGCGTTCGGCACAGGAGTACAGGATATATTTGCGCACT
GAAAGTAATGGAGAAGGAAGAAATAATAAAGTATAATTTACAGAAACAATTCAGAAGGGAGGTAGAAATA
CAAACATCGCTAAATCATCCGAATCTAACTAAATCATACGGCTATTTTCATGATGAAAAAAGAGTGTACC
TGCTAATGGAATACTTAGTCAATGGGGAAATGTATAAACTATTGAGGTTACACGGACCTTCAACGATAT
TTTAGCATCAGATTATATTTATCAAATTGCCAATGCCCTAGATTATATGCATAAAAAGAATATTATTCAT
AGAGATATTAACCTGAAAATATACTAATAGGGTTCAATAATGTCATTAAGTTAACGGACTTCGGATGGA
GTATAATAAATCCGCCAGAAAATAGAAGGAAAAGTGTCTGTGGGACAATTGACTACCTTTCTCCAGAAAT
GGTGGAGTCAAGGGAATATGATCACACTATAGATGCATGGGCTCTTGGCGTCCTGGCGTTTGAACACTG
ACCGGTGCCCTCCGTTTCGAAGAAGAAATGAAAGATACTACATATAAAAGGATAGCAGCACTGGATATCA
AATGCCCAGTAACATTTCTCAGGATGCGCAAGATTTAATACTTAACTACTAAAATACGACCCCAAAGA
TAGAATGCGCCTTGAGACGTAAAAATGCATCCTTGATACTAAGAAACAAGCCCTTTTGGGAAAATAAG
CGTTATAGAATTAAGTATGACAGAATCGTTTGAAGGGCACTATTAATCACTCCCGCACATATCACATA
ATACTAAGTATCCATTTCTAATATTTCACTCTTTTCGGCATCGTATATTGCGATATTTGATTAAATT
TTCTTGTTCAATTTTTCTCTTTTCTTTTCGCTTTGTGCGAAAGAAAAGAGGAAAACAAGCTGAAAATTGC
TATGCATTAAGTAGCAGATTTACTTTGTTGAGTTGGTTCTGATCAATAATAAGAGTAATGAAAGAAAGC
AAAAAATGGCTAAAGATAATTTAACTAATTTGCTCTCTCAATTGAACATTCAATTGTCTCAA
```

# The National Center for Biotechnology Information

- Created as a part of NLM in 1988
  - Establish public databases
  - Perform research in computational biology
  - Develop software tools for sequence analysis
  - Disseminate biomedical information

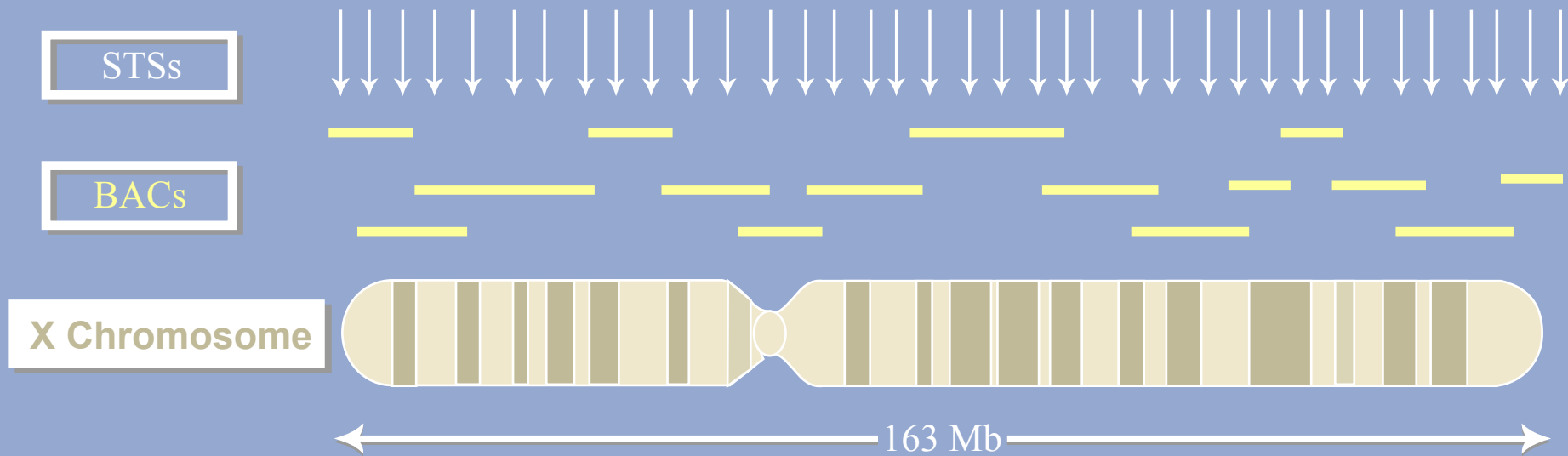
# Molecular Databases

- Primary Databases
  - Original submissions by experimentalists
  - Database staff organize but don't add additional information
    - **Example: GenBank**
- Derivative Databases
  - Human curated
    - compilation and correction of data
    - **Example: SWISS-PROT, NCBI RefSeq mRNA**
  - Computationally Derived
    - **Example: UniGene**
  - Combinations
    - **Example: NCBI Genome Assembly**

# What is GenBank?

## NCBI's Primary Sequence Database

- **Nucleotide only sequence database**
- **Archival in nature**
- **GenBank Data**
  - Direct submissions individual records (BankIt, Sequin)
  - Batch submissions via email (EST, GSS, STS)
  - ftp accounts sequencing centers
- **Data shared three collaborating databases**
  - GenBank
  - DNA Database of Japan (DDBJ).
  - European Molecular Biology Laboratory Database (EMBL) at EBI.



**Relationships of chromosomes to genome sequencing markers.**

The X chromosome is about 163 Mb in length. In this diagram, there are 16 overlapping BAC clones that span the entire length. In reality, 1,408 BACs were needed to span the X chromosome. Arrows (top) mark STSs scattered throughout the chromosome and on overlapping BACs.

# GenBank: NCBI's Primary Sequence Database

| <b>Release 133</b> | <b>December 2003</b> |
|--------------------|----------------------|
| 19,808,101         | Records              |
| 28,507,990,166     | Nucleotides          |
| 110,000 +          | Species              |

- full release every two months
- incremental and cumulative updates daily
- available only through internet

<ftp://ftp.ncbi.nih.gov/genbank/>

107.07 Gigabytes of data

# GenBank Divisions

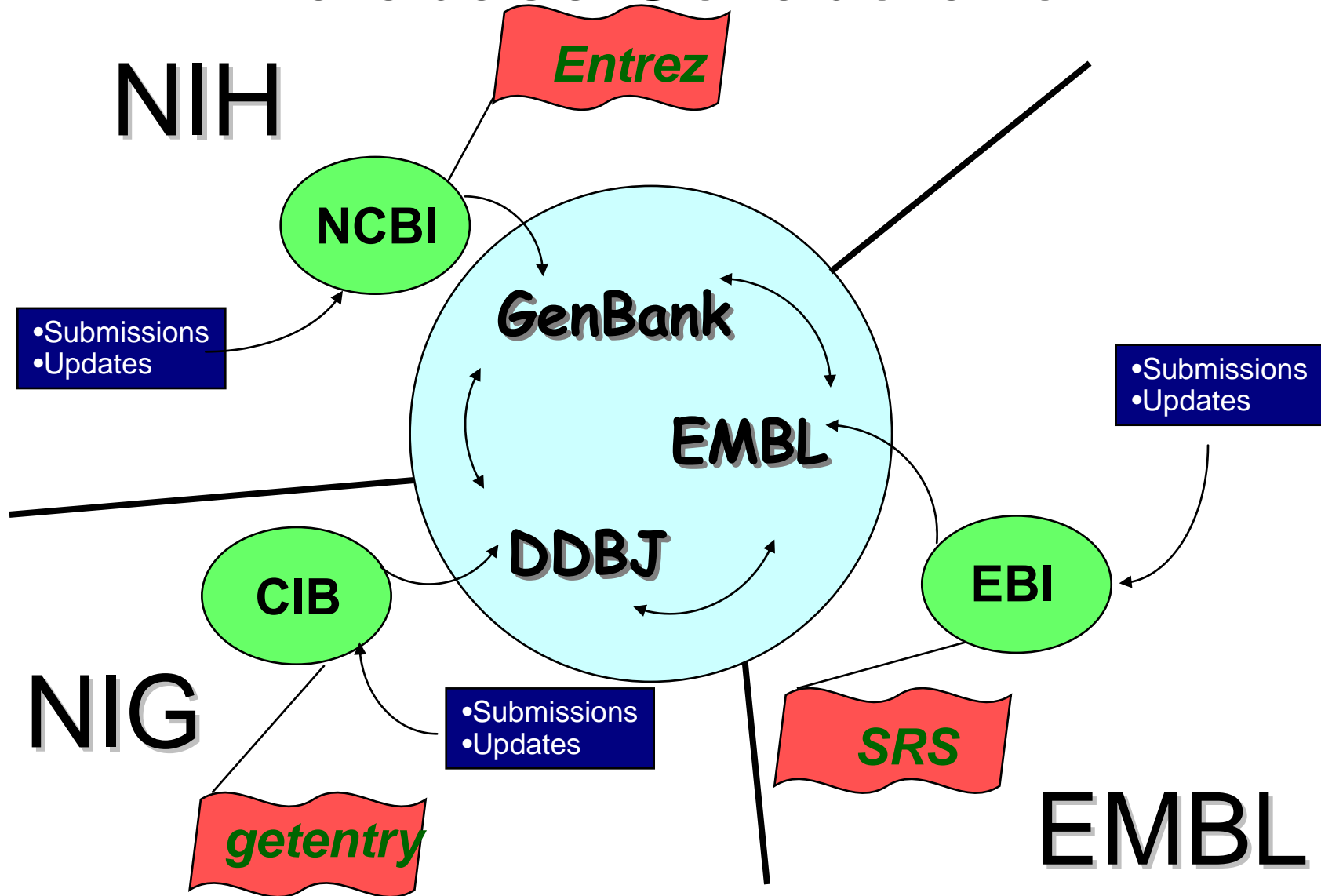
## Bulk Sequence Divisions

|            |                         |
|------------|-------------------------|
| <b>PAT</b> | Patent                  |
| <b>EST</b> | Expressed Sequence Tags |
| <b>STS</b> | Sequence Tagged Sites   |
| <b>GSS</b> | Genome Survey Sequences |
| <b>HTG</b> | High Throughput Genome  |
| <b>HTC</b> | High Throughput cDNA    |
| <b>CON</b> | Contig                  |

## Traditional Divisions

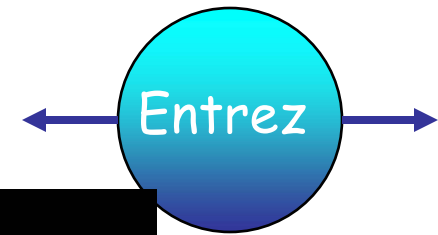
**BCT INV MAM PHG PLN PRI**  
**ROD SYN UNA VRL VRT**

# The International Sequence Database Collaboration





# Web Access Text



<http://www.ncbi.nlm.nih.gov>

NCBI Entrez, The Life Sciences Search Engine

HOME SEARCH SITE

|   |  |
|---|--|
| <b>PubMed:</b> biomedical literature citations and abstracts        | <b>Books:</b> online books                       |
| <b>PubMed Central:</b> free, full text journal articles             | <b>OMIM:</b> Online Mendelian Inheritance in Man |
| <b>Journals:</b> detailed information about journals in Entrez      | <b>Site Search:</b> NCBI web and FTP sites       |
| <b>MeSH:</b> detailed information about NLM's controlled vocabulary |  |

|   |  |
|---|--|
| <b>Nucleotide:</b> sequence database (GenBank)                | <b>UniGene:</b> gene-oriented clusters of transcript sequences |
| <b>Protein:</b> sequence database                             | <b>CDD:</b> conserved protein domain database                  |
| <b>Genome:</b> whole genome sequences                         | <b>3D Domains:</b> domains from Entrez Structure               |
| <b>Structure:</b> three-dimensional macromolecular structures | <b>UniSTS:</b> markers and mapping data                        |
| <b>Taxonomy:</b> organisms in GenBank                         | <b>PopSet:</b> population study data sets                      |
| <b>SNP:</b> single nucleotide polymorphism                    | <b>GEO:</b> expression and molecular abundance profiles        |
| <b>Gene:</b> gene-centered information                        | <b>GEO DataSets:</b> experimental sets of GEO data             |

Enter terms and click 'GO' to run the search against ALL the databases, OR  
Click Database Name or Icon to go directly to the Search Page for that database, OR  
Click Question Mark for a short explanation of that database.

[Disclaimer](#) | [Privacy statement](#) | [Accessibility](#)

## Sequence



## Structure



# Genetics experiment:

Isolate a yeast mutant that has increased chromosome number

phenotype

Can rescue the phenotype with a piece of DNA

...corresponds to a site of mutation in the yeasts DNA

genotype

```
CGTTTTCTGTAAAGCGCTTAATTTGTTTACCATTCTATAAAAACCTTGAGCTAAGGCCAACTGATGCA
ATTGCTCAAGTGAATGCATAAACAAAGCAAGATCATTCTTAGCGCAAAAAAACTGGGATTTTGAAATAC
AACAAAAGAAAGAAGTAAAAAGGGAATGCAACGCAATAGTTTAGTAAATATCAAACCTAAACGCTAATTCG
CCATCGAAAAAGACCACAACAAGACCAAATACGTCCAGGATCAATAAACCATGGAGAATATCCCATTCCG
CGCAGCAAAGAAACCCGAATTCAAAAATACCTTCACCTGTAAGAGAAAAATTGAACAGATTACCTGTAAA
CAATAAGAAGTTTTTGGATATGGAAAGCTCCAAAATTCCATCACCTATAAGGAAAGCGACTTCTTCCAAA
ATGATACACGAAAATAAGAAGCTACCTAAATTTAAATCCCTATCACTCGATGACTTTGAACTGGGGAAGA
AATTAGGAAAGGGTAAATTCGGTAAAGTTTATTGCGTTCGGCACAGGAGTACAGGATATATTTGCGCACT
GAAAGTAATGGAGAAGGAAGAAATAATAAAGTATAATTTACAGAAACAATTCAGAAGGGAGGTAGAAATA
CAAACATCGCTAAATCATCCGAATCTAACTAAATCATACGGCTATTTTCATGATGAAAAAAGAGTGTACC
TGCTAATGGAATACTTAGTCAATGGGGAAATGTATAAACTATTGAGGTTACACGGACCTTCAACGATAT
TTTAGCATCAGATTATATTTATCAAATTGCCAATGCCCTAGATTATATGCATAAAAAGAATATTATTCAT
AGAGATATTAACCTGAAAATATACTAATAGGGTTCAATAATGTCATTAAGTTAACGGACTTCGGATGGA
GTATAATAAATCCGCCAGAAAATAGAAGGAAAACCTGTCTGTGGGACAATTGACTACCTTTCTCCAGAAAT
GGTGGAGTCAAGGGAATATGATCACACTATAGATGCATGGGCTCTTGGCGTCCTGGCGTTTGAACACTG
ACCGGTGCCCTCCGTTTCGAAGAAGAAATGAAAGATACTACATATAAAAGGATAGCAGCACTGGATATCA
AATGCCCAGTAACATTTCTCAGGATGCGCAAGATTTAATACTTAACTACTAAAATACGACCCCAAAGA
TAGAATGCGCCTTGAGACGTAAAAATGCATCCTTGATACTAAGAAACAAGCCCTTTTGGGAAAATAAG
CGTTATAGAATTAAGTATGACAGAATCGTTTGAAGGGCACTATTAATCACTCCCGCACATATCACATA
ATACTAAGTATCCATTTCTAATATTTCACTCTTTTCGGCATCGTATATTGCGATATTTGATTAAATT
TTCTTGTTCAATTTTTCTCTTTTCTTTTCGCTTTGTGCGAAAGAAAAGAGGAAAACAAGCTGAAAATTGC
TATGCATTAAGTAGCAGATTTACTTTGTTGAGTTGGTTCTGATCAATAATAAGAGTAATGAAAGAAAGC
AAAAAATGGCTAAAGATAATTTAACTAATTTGCTCTCTCAATTGAACATTCAATTGTCTCAA
```

**NEW** 2 February 2004 BLAST 2.2.7 has been released. [Read more...](#)

Info

- FAQs
- News
- References
- Credits

Education

- Program selection guide
- Tutorial
- URL API guide

Download

- Executables
- Databases
- Source code

Support

- Helpdesk
- Mailing list

## Nucleotide

- Discontiguous megablast
- Megablast
- **Nucleotide-nucleotide BLAST (blastn)**
- Search for short nearly exact matches
- Search trace archives with megablast or discontiguous megablast

## Protein

- Protein-protein BLAST (blastp)
- PHI- and PSI-BLAST
- Search for short nearly exact matches
- Search the conserved domain database (rpsblast)
- Search by domain architecture (cdart)

## Translated

- Translated query vs. protein database (blastx)
- Protein query vs. translated database (tblastn)
- Translated query vs. translated database (tblastx)

## Genomes

- Human, mouse, rat
- Fugu rubripes, zebrafish
- Insects, nematodes, plants, lungi, malaria
- Microbial genomes, other eukaryotic genomes

## Special

- Align two sequences (bl2seq)
- Screen for vector contamination (VecScreen)
- Immunoglobulin BLAST (IgBlast)

## Meta

- Retrieve results by RID
- Get this page with javascript-free links

Search

```
AGGGTAAATTCGGTAAAGTTATTGCGTTCGGCACAGGAGTACAGGATATATTTCCGCACT  
CGAGAACGAAGAATAATAAAGTATAATTTACAGAAACAATTCAGAAGGGAGGTAGAATA  
CTAAATCATCCGAATCTAAC TAAATCATACGGCTATTTTCATGATGAAAAAAGAGTGTACC  
ATACTTACTCAATGGGGAAATCTATAAACTATTGAGGTTACACGGACCCTTCAACGATAT  
AGATTATATTTATCAAATTGCCAATGCCCTAGATTATATGCATAAAAAAGAATATTATCAT
```

Set subsequence From: \_\_\_\_\_ To: \_\_\_\_\_

Choose database  ▼

Now: **BLAST!** or **FASTA** **FASTX**

### Options for advanced blasting

Limit by entrez query \_\_\_\_\_ or select from:  ▼

Choose filter  Low complexity  Human repeats  Mask for lookup table only  Mask lower case

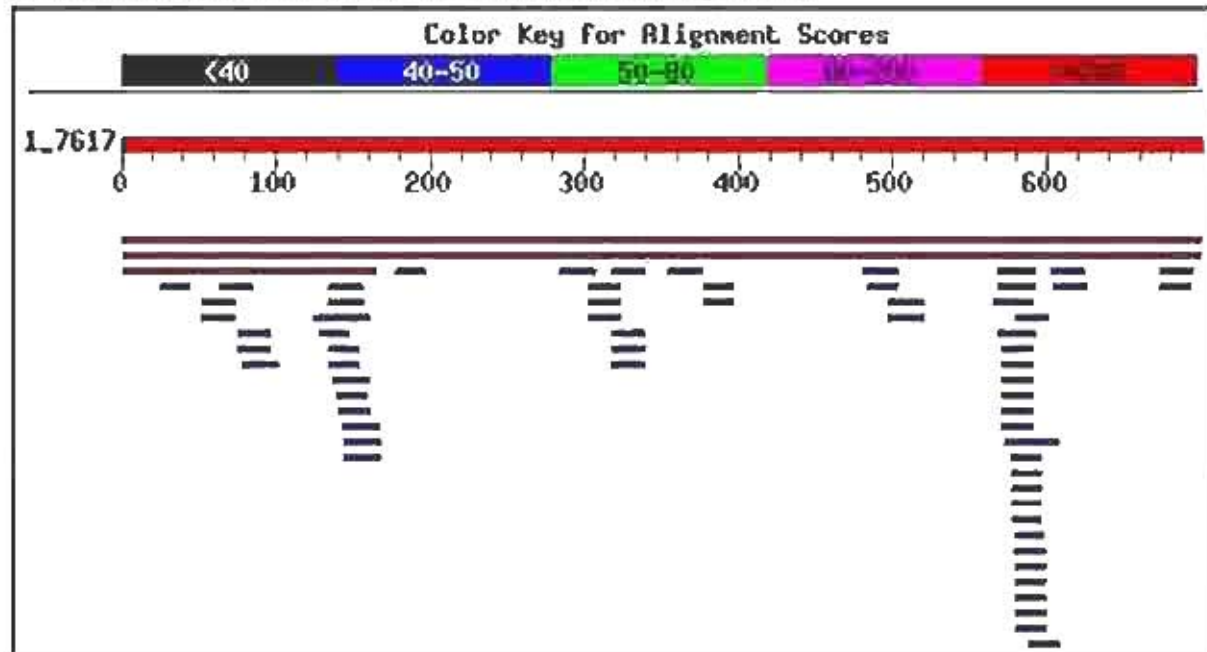
Evalue

Word Size  ▼

Other advanced \_\_\_\_\_

## Distribution of 67 Blast Hits on the Query Sequence

Mouse-over to show defline and scores. Click to show alignments



Sequences producing significant alignments:

|                                   |                                     |
|-----------------------------------|-------------------------------------|
| gi 1370433 emb Z73565.1 SCYPL209K | S.cerevisiae chromosome X...        |
| gi 560243 gb U07163.1 SCU07163    | Saccharomyces cerevisiae S28...     |
| gi 1370431 emb Z73564.1 SCYPL208K | S.cerevisiae chromosome X...        |
| gi 14164523 dbj AF003143.2        | Oryza sativa (japonica cultivar...  |
| gi 10257386 dbj AF002869.1        | Oryza sativa (japonica cultivar...  |
| gi 22296697 gb AC091043.12        | Homo sapiens chromosome 18, clo...  |
| gi 27356774 gb AC115119.5         | Mus musculus BAC clone RP23-B7M1... |
| gi 3900845 gb AC005078.1 AC005078 | Homo sapiens BAC clone CT...        |
| gi 37620357 gb AC146038.2         | Pan troglodytes chromosome 7 clo... |
| gi 24158583 gb AC121582.3         | Mus musculus BAC clone RP23-257I... |
| gi 38524392 emb BX004673.10       | Zebrafish DNA sequence from cl...   |

| Score (bits) | E Value |
|--------------|---------|
| 1340         | 0.0     |
| 1340         | 0.0     |
| 280          | 3e-72   |
| 46           | 0.089   |
| 46           | 0.089   |
| 46           | 0.089   |
| 43           | 0.35    |
| 43           | 0.35    |
| 43           | 0.35    |
| 42           | 1.4     |
| 42           | 1.4     |





[Get selected sequences](#) | [Select all](#) | [Deselect all](#)

Database: All GenBank+EMBL+DBJ+PDB sequences (but no EST, STS, GSS, or phase 0, 1 or 2 HTGS sequences)

Posted date: Jan 28, 2004 9:53 PM

Number of letters in database: 25,056,781

Number of sequences in database: 1,740,955

| Lambda | K     | H    |
|--------|-------|------|
| 1.37   | 0.711 | 1.11 |

| Gapped<br>Lambda | K     | H    |
|------------------|-------|------|
| 1.37             | 0.711 | 1.11 |

Matrix: blastn matrix:1 -1  
Gap Penalties: Existence: 5, Extension: 2

Number of hits to db: 242520

Number of Sequences: 2032434

Number of extensions: 242520

Number of successful extensions: 5340

Number of sequences better than 10.0: 1

Number of HSP's better than 10.0 without gapping: 3

Number of HSP's successfully gapped in prelim test: 0

Number of HSP's that attempted gapping in prelim test: 5330

Number of HSP's gapped (non-prelim): 10

length of query: 1402

length of database: 9,816,170,500

effective HSP length: 22

effective length of query: 678

effective length of database: 9,791,456,952

effective search space: 6638607813456

effective search space used: 6638607813456

T: 0

A: 0

X1: 6 (11.9 bits)

X2: 15 (29.7 bits)

S1: 12 (24.3 bits)

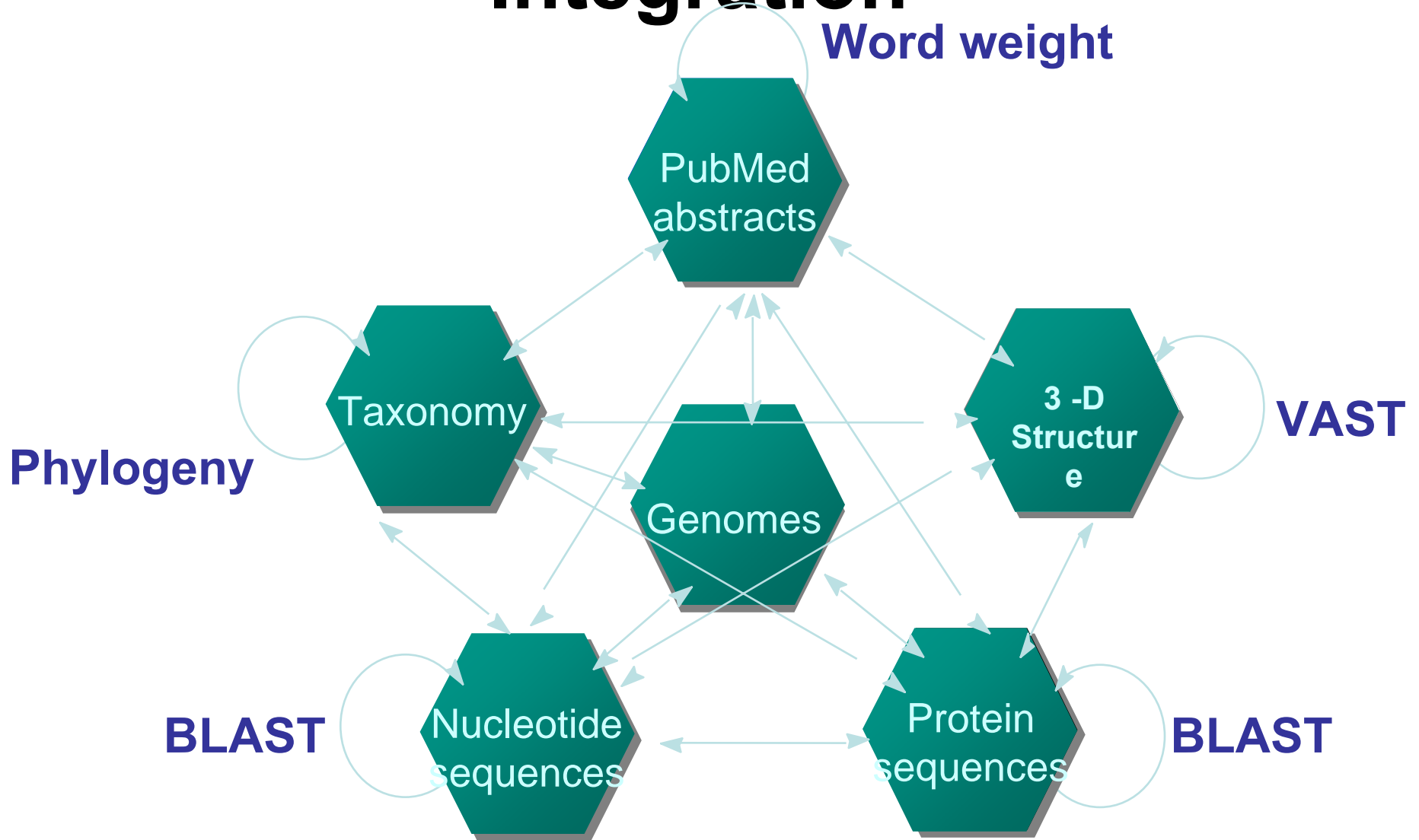
S2: 20 (40.1 bits)

# Using Entrez

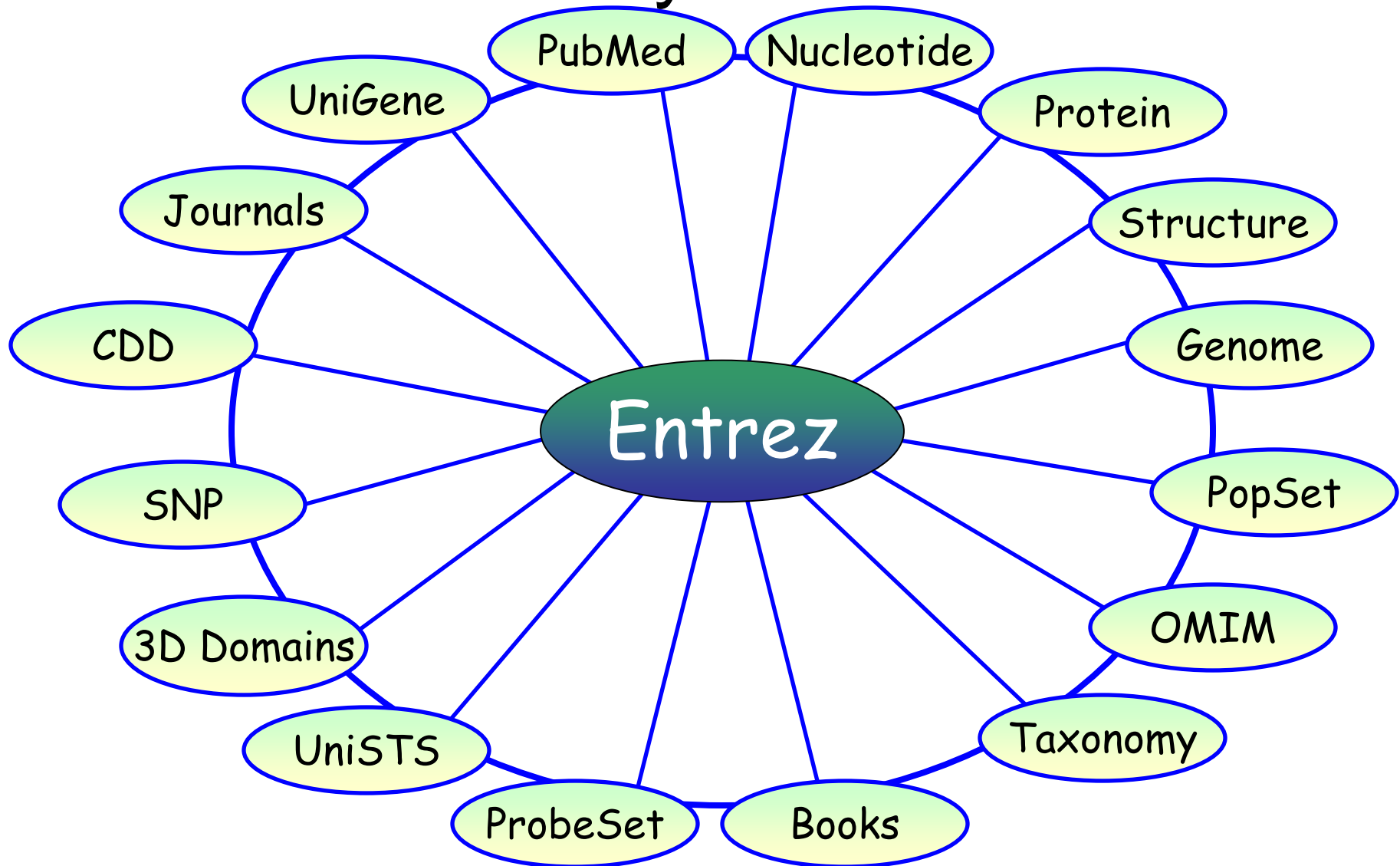
An integrated database  
search and retrieval system



# Entrez: Database Integration



# The (ever) Expanding Entrez System



# Entrez Databases

|            |   |
|------------|---|
| PubMed     | Biomedical literature                                       |
| Books      | Online textbooks  |
| Nucleotide | GenBank, EMBL, DDBJ, RefSeq, PDB                            |
| Protein    | [GenBank, EMBL, DDBJ], RefSeq,<br>SWISS-PROT, PIR, PRF, PDB |
| Genome     | Complete genomes  |
| Taxonomy   | Organisms in NCBI sequence databases                        |
| Structure  | MMDB: experimental 3D structures                            |
| Domains    | CDD: conserved protein domains                              |
| 3D Domains | Compact 3D protein domains in MMDB                          |
| OMIM       | Online Mendelian Inheritance in Man                         |
| SNP        | Single nucleotide polymorphisms                             |
| UniSTS     | Sequence Tagged Site markers                                |
| ProbeSet   | Gene expression and microarray datasets                     |
| PopSet     | Population study datasets                                   |
| UniGene    | Gene-based expressed sequence clusters                      |

# A Traditional GenBank Record

Back Forward Mail

NCBI Sequence Viewer LocusLink UniGene Entrez GEO - Gene Expressi...

NCBI

Entrez PubMed Nucleotide Protein Genome Structure PMC Taxonomy Books

Search Nucleotide for

Limits Preview/Index

Display default Show: 20 Send to File Get

**Definition = Title**

I: [U07163](#). *Saccharomyces cerevisiae* [gi:460243]

LOCUS SCU07163 1603 bp DNA linear PLN 04-AUG-1994

DEFINITION *Saccharomyces cerevisiae* S288C Ipl1p protein kinase (IPL1) gene, complete cds.

ACCESSION U07163

VERSION U07163.1 GI:460243

KEYWORDS .

SOURCE *Saccharomyces cerevisiae* (baker's yeast)

ORGANISM [Saccharomyces cerevisiae](#)  
Eukaryota; Fungi; Ascomycota; Saccharomycotina; Saccharomycetes; Saccharomycetales; Saccharomycetaceae; Saccharomyces.

REFERENCE 1 (bases 1 to 1603)

|           |                    |
|-----------|--------------------|
| ACCESSION | U07163             |
| VERSION   | U07163.1 GI:460243 |

**Accession Number**

**Version Number**

**GI Number**

**NCBI's Taxonomy**

AUTHOR Chan, C.S.

TITLE Direct Submission

JOURNAL Submitted (25 FEB 1994) Lawrence S.M., Department of Experimental Science, Texas at Austin, TX 78712

FEATURES

source

organism="Saccharomyces cerevisiae"

mol\_type="genomic DNA"

strain="S288C"

db\_xref="taxon:4932"

clone="pCC36"

clone\_lib="YCP50 library"

166..1269

gene="IPL1"

166..1269

gene="IPL1"

function="yeast c"

codon\_start=1

product="Ipl1p protein kinase"

protein\_id="AAA20496.1"

db\_xref="GI:460244"

gene

CDS

# FASTA Format

Display FASTA Show: 20 Send to File Get Subsequent

1: **FASTA Definition Line**

>gi|460243|gb|U07163.1|SCU07163

gi number

Accession number

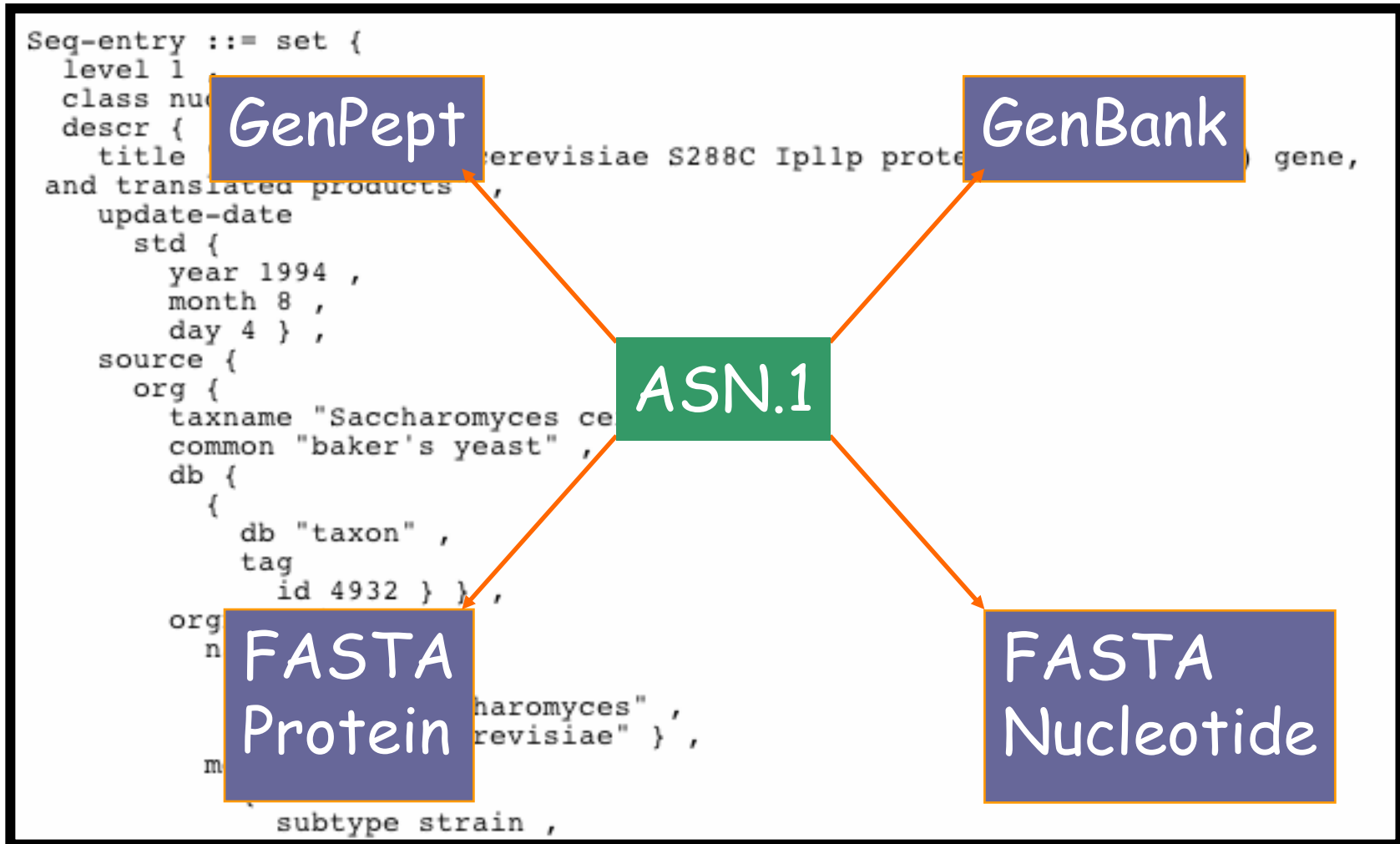
Locus Name

**Database Identifiers**

|     |                  |
|-----|------------------|
| gb  | GenBank          |
| emb | EMBL             |
| dbj | DDBJ             |
| sp  | SWISS-PROT       |
| pdb | Protein Databank |
| pir | PIR              |
| prf | PRF              |
| ref | RefSeq           |

CGTTT  
ATTGC  
AACAA  
CCATCGAAAAAGACAACAAGACC  
CGCAG  
CAATA  
ATGATACACGAAATAGAA  
AATTAGGAAAGGGTAAATT  
GAAAGTAATGGAGAAGGAA  
CAAACATCGCTAAATCATC  
TGCTAATGGAATACTTAGT  
TTTAGCATCAGATTATATT  
AGAGATATTAACCTGAAA  
GTATAATAAATCCGCCAGA  
GGTGGAGTCAAGGGAATAT  
ACCGGTGCCCTCCGTTCCG  
AAATGCCCAGTAACATTTCT  
TAGAATGCGCCTTGGAGAC  
CGGTTATAGAATTAAGTA  
ATACTAAGTATCCATTTCT  
TTCTTGTTCATTTTTCT  
TATGCATTAAAGTAGCAGA  
AAAAAATGGCTAAAGATA

# Abstract Syntax Notation: ASN.1



# NCBI Toolbox

```
/*
 *
 *   asn2ff.c
 *   convert an ASN.1 entry to flat file format, using the FFPrintArray.
 *
 ****
#include <accent
#include "asn2ff
#include "asn2ff
#include "ffprin
#include <subuti
#include <objall
#include <objcod
#include <lsqfet
#include <explor

#ifdef ENABLE_ID1
#include <accid1
#endif

FILE *fpl;

Args myargs[] = {
  "Input asnfile in binary mode", "F", NULL, NULL, TRUE, 'b', ARG_BOOLEAN, 0.0, 0, NULL},
  "Output Filename", "stdout", NULL, NULL, TRUE, 'o', ARG_FILE_OUT, 0.0, 0, NULL},
  "Show Sequence?", "T", NULL, NULL, TRUE, 'h', ARG_BOOLEAN, 0.0, 0, NULL},

```

## Toolbox Sources

```
ftp> open ftp.ncbi.nih.gov
```

.

.

```
ftp> cd toolbox
```

```
ftp> cd ncbi_tools
```

```
ftp://ftp.ncbi.nlm.gov/toolbox/ncbi_tools
```

[gene](#)

[CDS](#)

```

/clone="pCC36"
/clone_lib="YcP50 library of Rose et al."
166..1269
/gene="IPL1"
166..1269
/gene="IPL1"
/function="yeast chromosome segregation"
/codon_start=1
/product="Ipl1p protein kinase"
/protein_id="AAA20496.1"
/db_xref="GI:460244"
/translation="MQRNSLVNIKLNANSPSKKTTTRPNTSRINKPWRISHSPQQRNP
NSKIPSPVREKLNRLPVNNKKFLDMESSKIPSPIRKATSSKMIHENKKLPKFKSLSD
DFELGKKLGKGFVKVYCVRHRSTGYICALKVMKEEIIKYNLQKQFRREVEIQTSLN
HPNLTKSYGYFHDEKRVYLLMEYLVNGEMYKLLRLHGPFDNILASDYIQIANALDYM
HKKNIHRDIKPENILIGFNNVIKLTDFGWSIINPPENRRKTVCGTIDYLSPEMVESR
IDIKMPSNISQDAQDLILK

```

**/protein\_id="AAA20496.1"**

**/db\_xref="GI:460244"**

**GenPept Protein IDs**

```

ttg agtaaggee
61 aactgatgca attgctcaag tgaatgcata aacaaagcaa gatcattcct agcgcaaaaa
121 aaactgggat ttgaaatac aacaaaagaa agaagtaaaa agggaatgca acgcaatagt
181 ttagtaata tcaactaaa cgtaattcg ccatcgaaaa agaccacaac aagaccaaat
241 acgtccagga tcaataaac atggagaata tccatttcgc cgcagcaaaag aaaccggaat
301 tcaaaaatac cttcacctgt aagagaaaaa ttgaacagat tacctgtaaa caataagaag
361 tttttggata tggaaagctc caaaattcca tcacctataa ggaaagcgac ttcttccaaa
421 atgatacacg aaaataagaa gctacctaaa tttaaatccc tatcactcga tgactttgaa
481 ctgggggaaga aattaggaaa gggtaaatc ggtaagttt attgcgttcg gcacaggagt
541 acaggatata tttgcgcact gaaagtaatg gagaaggaag aaataataaa gtataattta
601 cagaaacaat tcagaaggga ggtagaata caaacatcgc taaatcatcc gaatcctaact
661 aatcatatcg gctattttca tgatgaaaaa agagtgtacc tgctaattgga atacttagtc
721 aatggggaaa tgtataaact attgaggtta cacggacct tcaacgatat tttagcatca
781 gattatattt atcaaatgac caatgcctca gattatatgc ataaaaagaa tattattcat

```





Protein

for

1: AAA20496 [Ipl1 protein kin. |g:460244]



LOCUS AAA20496 387 aa linear PLN 04-AUG-1994

DEFINITION Ipl1p protein kinase.

ACCESSION AAA20496

VERSION AAA20496.1 GI:460244

DESCRIPTION locus BCID7191 accession [BC07161](#)

KEYWORDS .

SOURCE Saccharomyces cerevisiae (baker's yeast)

ORGANISM [Saccharomyces cerevisiae](#)

Eukaryota; Fungi; Ascomycota; Saccharomycotina; Saccharomycetes;

Saccharomycetales; Saccharomycetaceae; Saccharomycetes.

REFERENCE 1 (residues 1 to 387)

AUTHORS Francisco, L., Wang, V. and Chan, C.S.

TITLE Type 1 protein phosphatase acts in opposition to Ipl1 protein

kinase in regulating yeast chromosome segregation

JOURNAL Mol. Cell. Biol. 14 (7), 4711-4740 (1994)

RECORDING [24273879](#)

DOI [1007278](#)

REFERENCE 2 (residues 1 to 387)

AUTHORS Chan, C.S.

TITLE SAKKO Submission

JOURNAL Submitted (25-FEB-1994) Clarence H.S. Chan, Department of



Address <http://www.ncbi.nlm.nih.gov/blast/blast.cgi?seq=AA034961>

**FASTA** [FASTA](#)  
**REFERENCE** 2 (residues 1 to 367)  
**AUTHORS** Chan, C. S.  
**TITLE** Direct Substitution  
**JOURNAL** Submitted (25-FEB-1994) Clarence S.R. Chan, Department of Microbiology, University of Texas at Austin, Experimental Science building, Room 236, Austin, TX 78712, USA  
**COMMENT** Method: conceptual translation.

**FEATURES** Location/Qualifiers  
 source 1..367  
 /organism="Saccharomyces cerevisiae"  
 /strain="SC08C"  
 /db\_xref="taxon:4932"  
 /clone="pC026"  
 /clone\_lab="YCp50 library of Rose et al."  
[Protein](#) 1..367  
 /product="Ipl1 protein kinase"  
 /function="yeast chromosome degradation"  
[CDS](#) 1..367  
 /gene="IPL1"  
 /code1\_by="007163.1:166..1269"

**ORIGIN**

```

1  MPTNGLVNAK  LNAAPPKK  TERPTARIN  KPTSLHPG  QPAPLAKIP  PVI-KINCIP
61  VNAKFLAME  SSKIPQIK  ALEKMLN  XKIPKFALE  LKIFLAKKI  SKKFKVYC
121  VHRSTGYIC  AIXMKEEL  IYNIQKQC  REVERQAIN  HYNIKQYQ  IHDKRYVIL
181  MRYLVNGMY  KLIKILPFA  DILASDYIQ  IANALDYMK  KNIRHDIK  ENLIGFARV
241  IKLIDFYVI  IAPPNRYK  VQYFDYIAP  MVEVREYDH  TIDAVLGI  AFILKQAPP
301  FVEAKDITY  KRLELDIK  PNTLQDQD  LIKLIKYS  KDRMLQDK  AWPILRDKP
361  IWAARL

```

//

Retrieved July 5, 1992

**NCBI**

|        |            |             |        |          |
|--------|------------|-------------|--------|----------|
| BLAST  | Protein    | Structure   | PubMed | Taxonomy |
| Genome | Nucleotide | 3D-Database | Books  | Help     |

Query: gn460244 Regulation of yeast chromosome segregation: *hlp1* [Saccharomyces cerevisiae]  
 Matching GI: 639201, 122955, 1370424, 6326047

COG0016 assigned by Cogiter (3 best hits)

200 BLAST hits to 51 unique species [Sort by taxonomy proximity](#)

Archaea
  Bacteria
  Eukarya
  Metazoa
  Fungi
  Plants
  Viruses
  Other Eukaryotes

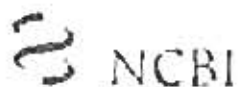
Keep only

367 aa

| Hit # | SCOPE | E | ACCESSION | GI       | PROTEIN DESCRIPTION                                  |
|-------|-------|---|-----------|----------|--|
| 1     | 639   | 4 | CAA10315  | 2995321  | ser/thr protein kinase [Schizosaccharomyces pombe]   |
| 2     | 727   | 4 | CAI78915  | 609282   | glo265 (Xenopus laevis)                              |
| 3     | 731   | 4 | CAI819    | 17921859 | serine/threonine protein kinase Eg2-like (p16X/Eg22) |
| 4     | 731   | 4 | CAA78919  | 609280   | p16X/Eg22 [Xenopus laevis]                           |
| 5     | 728   | 4 | AAM11289  | 12554871 | serine/threonine kinase 15 [Homo sapiens]            |
| 6     | 725   | 4 | AAC63902  | 1213197  | serine/threonine kinase [Homo sapiens]               |
| 7     | 726   | 4 | AAC12708  | 2078888  | aurore-related kinase 1 [Homo sapiens]               |
| 8     | 724   | 4 | BAA23592  | 2641040  | aurore/IPL1-related kinase [Homo sapiens]            |
| 9     | 721   | 4 | AAD78715  | 11866530 | aurore B [Xenopus laevis]                            |
| 10    | 723   | 4 | AAG10287  | 9686777  | protein kinase AIPK1 [Xenopus laevis]                |
| 11    | 721   | 4 | J02975    | 715611   | aurore-related kinase 1 (EC 2.7.11) - human          |







# NCBI Conserved Domain Summary

- [View Summary](#)
- [FASTA](#)
- [Nucleotide](#)
- [Protein](#)
- [Structure](#)
- [CDs](#)
- [Taxonomy](#)
- [Help](#)

**Query-** [U11360219:U11360219.1](#) Iplip protein kinase  
1347 amino acids

**Database:** cdd.v.1.60

Click on boxes for multiple alignments



[Links](#) | [Classifiers](#) | [Write in the Public Domain](#)  
[PubMed](#) | [PubMed](#) | [PubMed](#)



# CDD



PubMed

BLAST

OMM

Taxonomy

Entrez Structure

Search Entrez

Structure

for

Go

### CDD Help

### NCBI Handbook

Help on CD-Search and Database

### CD-Search

Search with advanced options

### CDART

Conserved Domain Architecture Retrieval Tool

### Smart

Protein Structure

## A Conserved Domain Database and Search Service, v1.60

Proteins often contain several modules or domains, each with a distinct evolutionary origin and function. The CD-Search service may be used to identify the conserved domains present in a protein sequence.

### Run CD-Search

Search Database:  Submit Query

Enter Protein Query as Accession, GI, or Sequence in FASTA format

Read about [FASTA format description](#), click [here](#) for advanced options

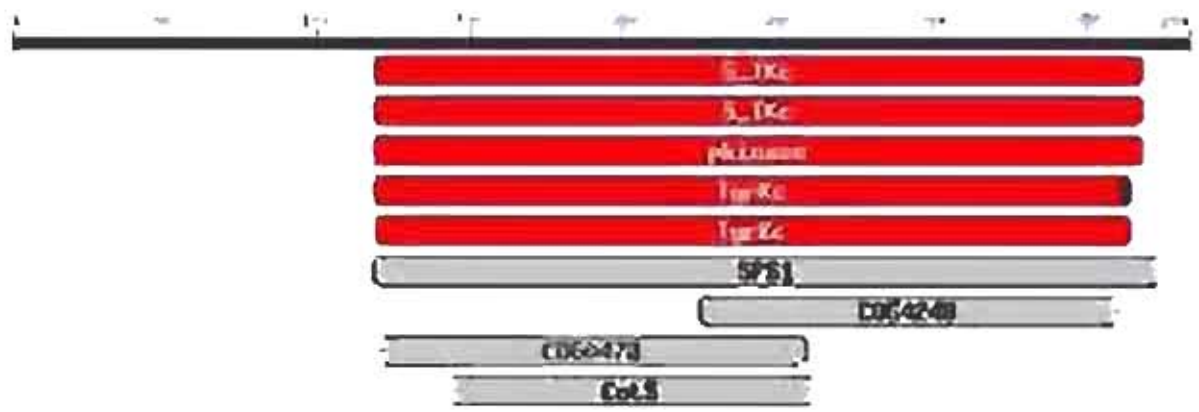
Computational biologists define conserved domains based on recurring sequence patterns or motifs. CDD currently contains domains derived from two methods:

Find CDS

by keyword



Click on boxes for multiple alignments



Show Domain Relatives

This CD alignment includes 3D structure. To display structure, download [Cn3D!](#)

| PSSMs producing significant alignments |   | Score  | E     |
|--|---|--------|-------|
|  |   | (bits) | value |
| <a href="#">cd0130</a>                 | S_TKc, Serine/Threonine protein kinases, catalytic do       | 304    | 8e-84 |
| <a href="#">smat00220</a>              | S_TKc, Serine/Threonine protein kinases, catalytic          | 302    | 5e-93 |
| <a href="#">pfam00669</a>              | pkinase, Protein kinase domain                              | 280    | 2e-70 |
| <a href="#">cd00192</a>                | TyrKc, Tyrosine kinase, catalytic domain, Phosphotran       | 160    | 4e-42 |
| <a href="#">smat00219</a>              | TyrKc, Tyrosine kinase, catalytic domain, Phosphot          | 160    | 6e-42 |
| <a href="#">COG515</a>                 | SPS1, SPS1, Serine/threonine protein kinase [General functi | 139    | 5e-39 |
| <a href="#">COG4248</a>                | COG4248, COG4248, Unclassified protein with protein kinase  | 93.2   | 5e-05 |



- [cd0190](#), smart00190, S\_TKc, Serine/Threonine protein kinases, catalytic domain, Phosphotransferases, Serine or threonine-specific kinase subfamily. The enzymatic activity of these protein kinases is controlled by phosphorylation of specific residues in the activation segment of the catalytic domain, sometimes combined with reversible conformational changes in the C-terminal autoregulatory tail.

CD-length = 257 residues, 99.1% aligned  
Score = 304 bits (781), Expect = 6e-64

```

Query: 110 F E I G Y P L G R G K T G R V T L L K I R E T Y F I V A L E T L H N S E L V G K I E V G V P P E I S I Q S L R S Y H 177
Subject: 2 Y E L L Q V L G G Q A F G R V T L A R D E R T Y E I V S T H I I K R E L S K - R Q V E P I L G E T K L K P I D N P H 60

Query: 170 I L P L T C H F H I R K P I V L I L Y A G G Q E L Y G M P S A K Y F H E I V A S Y I F L H A N A L S T L E R K H V 237
Subject: 61 I V K L T V V Y E R E D E L V L V H E Y C S G G S L F D L L K R P O P L S E D Y A P F Y F P Q I L S A L K V L R E D S I 120

Query: 230 I H E D I F P E N I L L G I S - E I I E L S G F Q H E - - V H A P D N N H T L C U T L D T L P P E H V T L H E H T E R V 295
Subject: 121 I H E D L A P E N I L L G E R H V L A D F V L A N C L O S G G T K L T F T N T P F L H A S E L L G Q Q Y G F A V 180

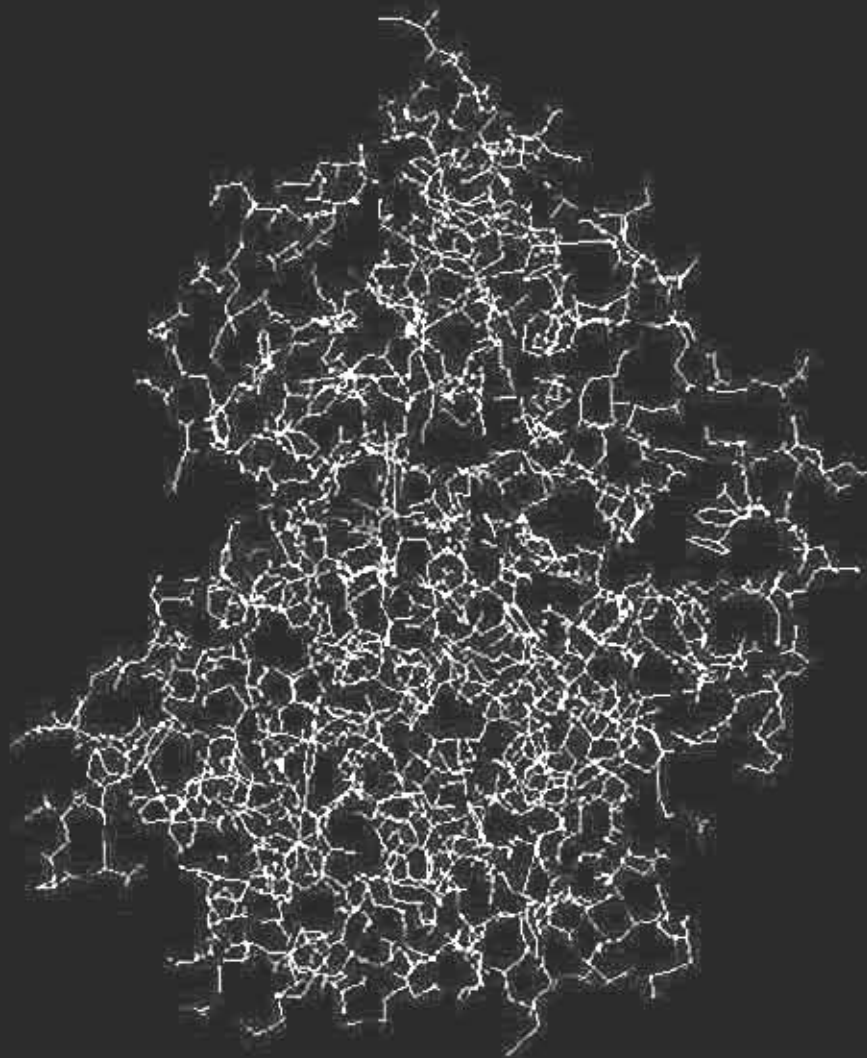
Query: 290 P L H E L G V L T E F I V S A P P Y C O M G C H A T Y H D L A F V O L E - - I H P V V P P E A P L I S U L L Q M 353
Subject: 181 S I H E L G V I L Y E L L Y T Y P P Y P G S N - E I E L L E K I L E G O L T P D E L F R K I Y P E A N D L I K L L A Y 230

Query: 353 N P E R M H L E Q V T R A P V I 349
Subject: 240 I P E R P L T A S Y A L N P S Y 206
  
```

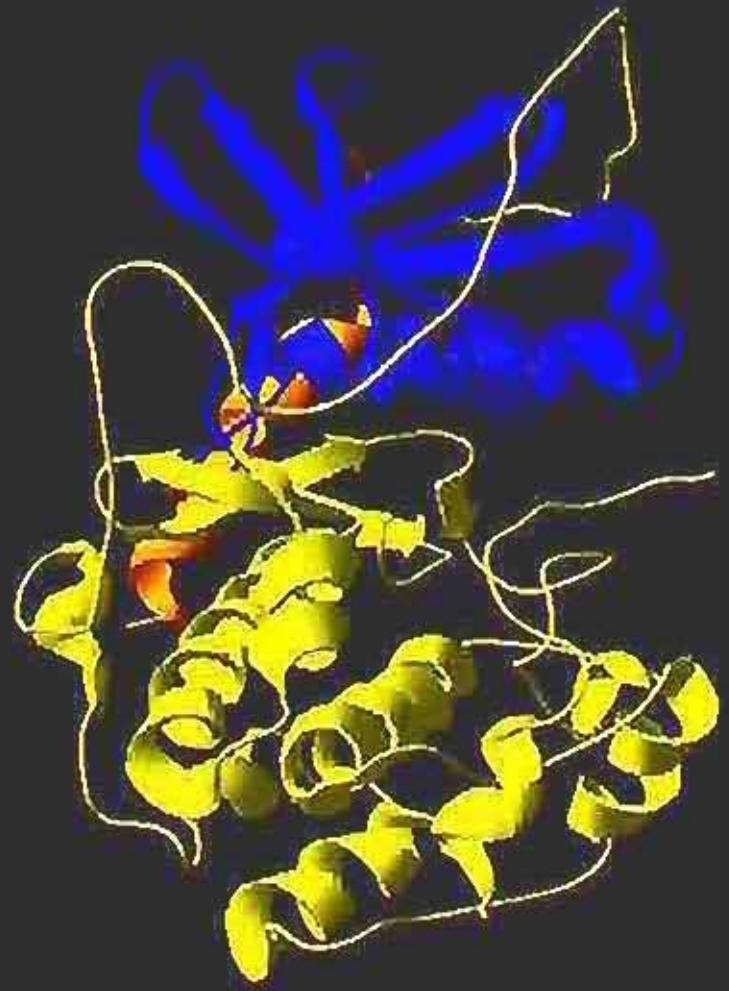
- [cd0220](#), smart00220, S\_TKc, Serine/Threonine protein kinases, catalytic domain, Phosphotransferases, Serine or threonine-specific kinase subfamily.

CD-length = 256 residues, 100.0% aligned  
Score = 302 bits (774), Expect = 1e-63









BREAK

# Why align sequences?

- Functional predictions based on identifying homologues.

Assumes:

conservation of  
sequence



conservation of  
function

**BUT:** Function carried out at level of proteins, i.e.  
3-D structure

Sequence conservation carried out at level of DNA  
1-D sequence


# Implicit Assumption of Evolution in Models of Sequence Homology

**Assume:**

Sequence conservation  Structure conservation

---

***Note that the converse is NOT necessarily true!***

Structure conservation  Sequence conservation

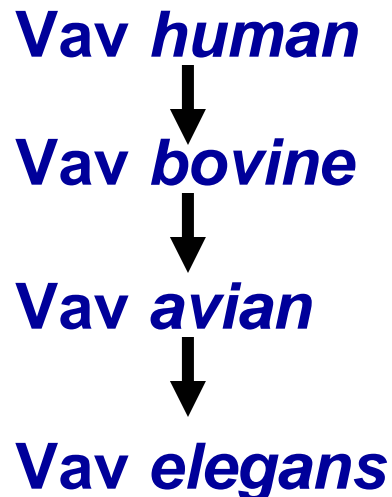
# Definitions

- Homologue:

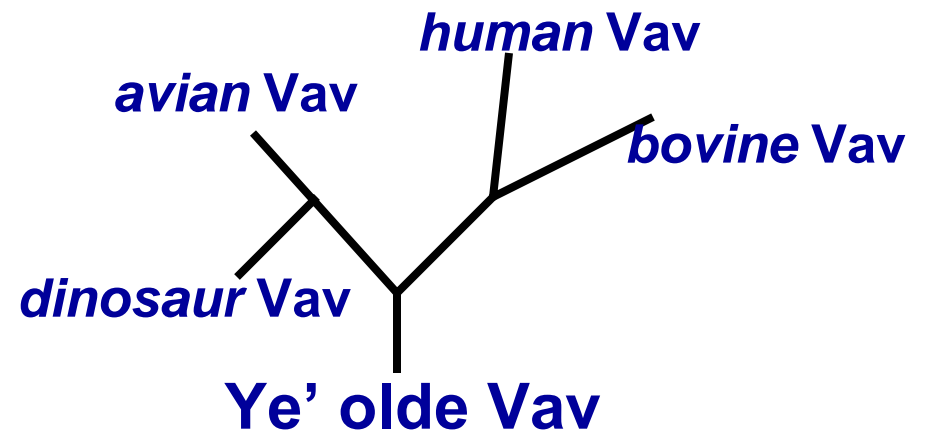
a relatively non-specific term (meaningless?) that conveys the idea that two sequences are somehow related

- Orthologue:

Ortho = (*greek*) straight....implies direct descent, 1 ancestor



**-or-**

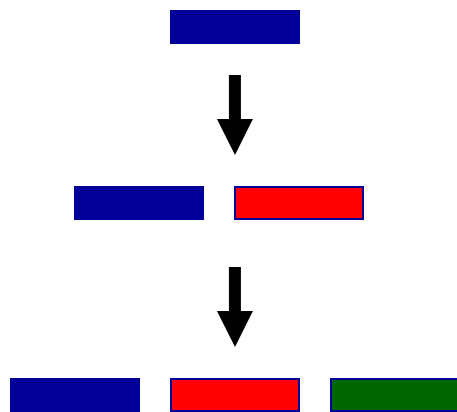


# Definitions

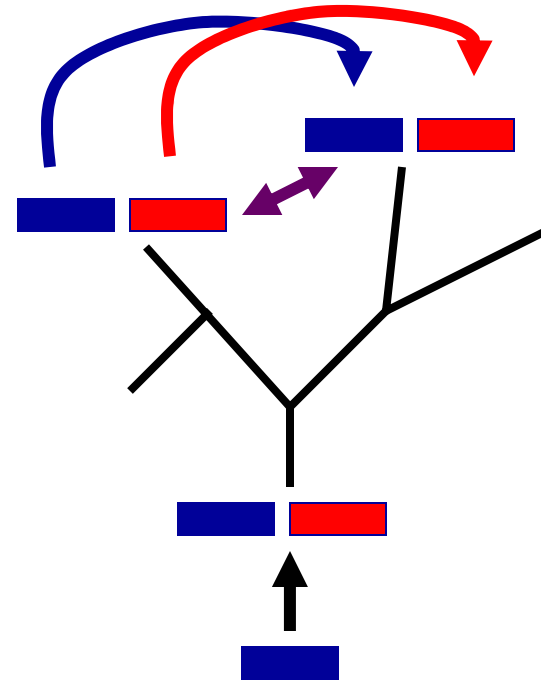
- Paralogue:

Para = (*greek*) along side of...implies some type of gene duplication event

*Vav1 human* ↔ *Vav2human* ↔ *Vav3human*

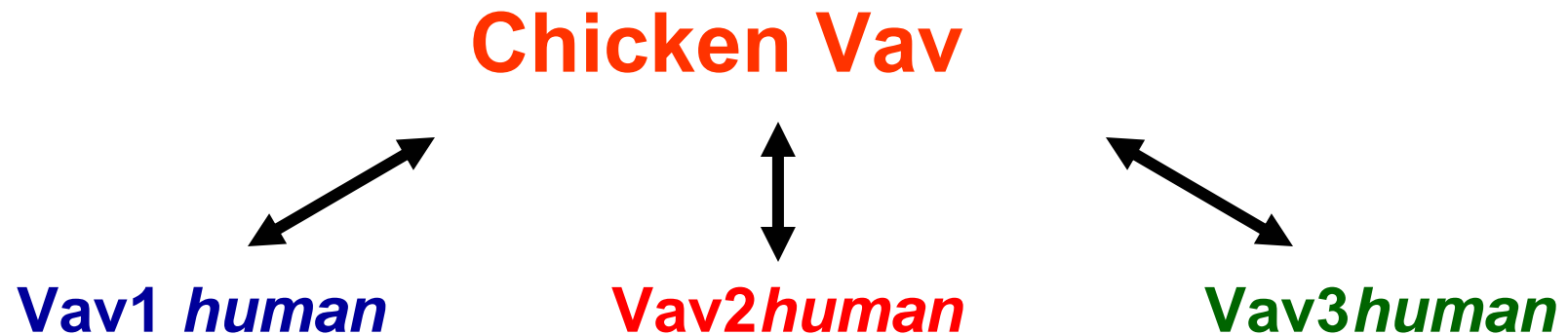


-or-





# Complex problem



*Orthologue or Parologue?*

**Need to know evolutionary relationships**

**Can try and get these from alignments as well...**

# Alignments- Good Bad and Ugly

*HgbA-human*

GSAQVKGHGKKVADALTNVAHAVDDMPNALSALSDLHAHKL  
: + + : : + : : : : + : + + + + : : + : + + + + + : : + : : + : : : : + : :  
GNPKVKAHGKKVLGAFSDGLAHL DNLKGT FATLSELHCDKL

*HgbB-human*

# Alignments- Good Bad and Ugly

## SPURIOUS ALIGNMENT

*HgbA-human*

GSAQVKGHGKKVADALTNVAHAVDDMPNALSALSD----LHAHKL  
::+ ++: + ++::: ++ :+ :+ : +++: +  
GSGYLVGDSLTFVDLL--VAQHTADLLAANAALLDEFPPQFKAHQE

*Nematode glutathione S-transferase*

*HgbA-human*

GSAQVKGHGKKVADALTNVAHV---D--DMPNALSALSDLHAHKL  
++ ++++:+ ::+ ++ +:++ + +: +:++ :  
NNPELQAHAGKVFKLVEAAIQLQVTGVVVTDATLKNLGSVHVS KG

*Leghemoglobin, yellow lupin*

# Alignments

- **Types:**

- **Local**
- **Global**
- **Ungapped**
- **Gapped (2 types- linear, affine)**

- **Methods:**

- **Dot matrix**
- **Dynamic Programming**
- **Word, k-tup**



## Example – self alignment

|   | G | F | D | S | F | K | R | L | E | F | S | E | V |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| G | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| F |   | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 |
| D |   |   | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| S |   |   |   | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 |
| F |   |   |   |   | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 |
| K |   |   |   |   |   | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| R |   |   |   |   |   |   | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| L |   |   |   |   |   |   |   | 1 | 0 | 0 | 0 | 0 | 0 |
| E |   |   |   |   |   |   |   |   | 1 | 0 | 0 | 1 | 0 |
| F |   |   |   |   |   |   |   |   |   | 1 | 0 | 0 | 0 |
| S |   |   |   |   |   |   |   |   |   |   | 1 | 0 | 0 |

# Example – self alignment....with sliding window

|   | G | F | D | S | F | K | R | L | E | F | S | E | V |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| G | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| F |   | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 |
| D |   |   | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| S |   |   |   | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 |
| F |   |   |   |   | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 |
| K |   |   |   |   |   | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| R |   |   |   |   |   |   | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| L |   |   |   |   |   |   |   | 1 | 0 | 0 | 0 | 0 | 0 |
| E |   |   |   |   |   |   |   |   | 1 | 0 | 0 | 1 | 0 |
| F |   |   |   |   |   |   |   |   |   | 1 | 0 | 0 | 0 |
| S |   |   |   |   |   |   |   |   |   |   | 1 | 0 | 0 |

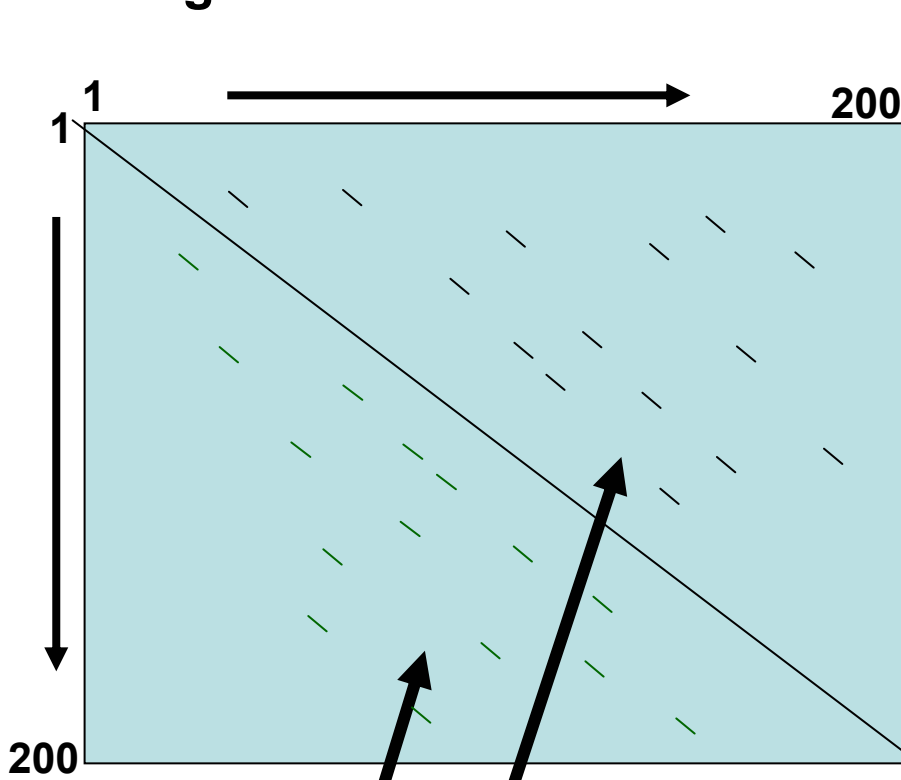
# Example – self alignment

|   | G | F | D | S | F | K | R | L | E | F | S | E | V |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| G | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| F |   | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| D |   |   | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| S |   |   |   | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| F |   |   |   |   | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| K |   |   |   |   |   | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| R |   |   |   |   |   |   | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| L |   |   |   |   |   |   |   | 1 | 0 | 0 | 0 | 0 | 0 |
| E |   |   |   |   |   |   |   |   | 1 | 0 | 0 | 0 | 0 |
| F |   |   |   |   |   |   |   |   |   | 1 | 0 | 0 | 0 |
| S |   |   |   |   |   |   |   |   |   |   | 1 | 0 | 0 |

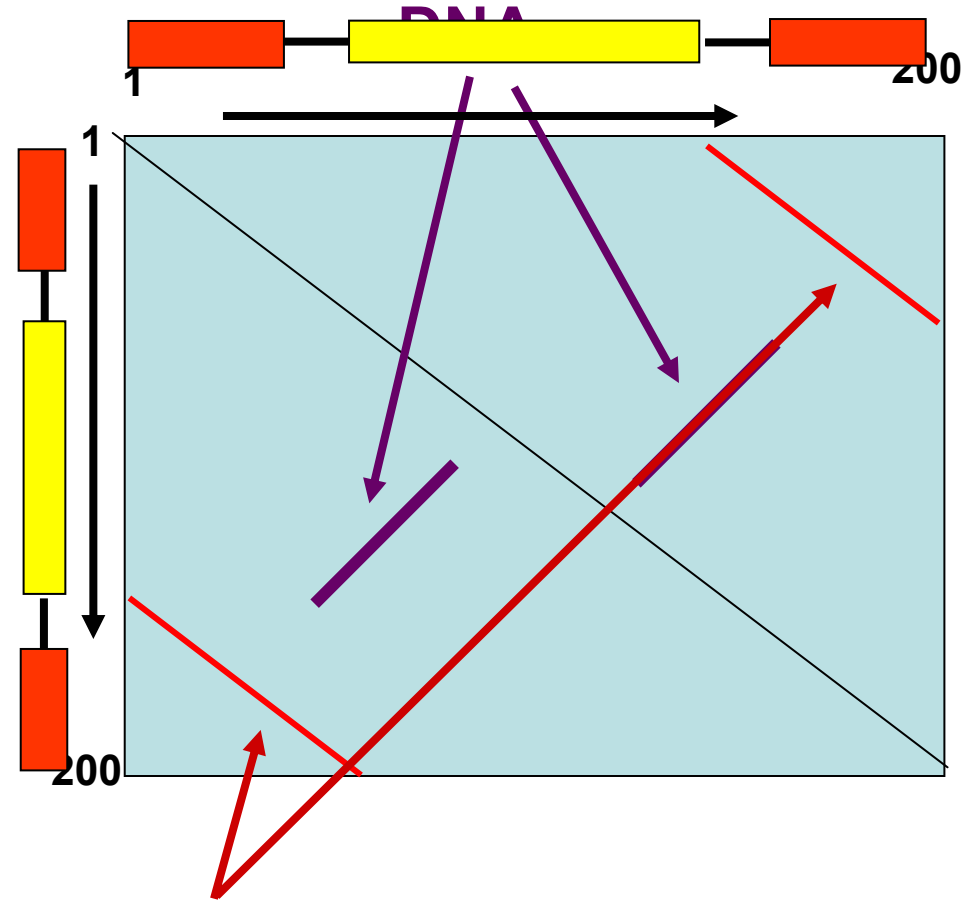


# Simple alignments

Self alignment – Dot Matrix



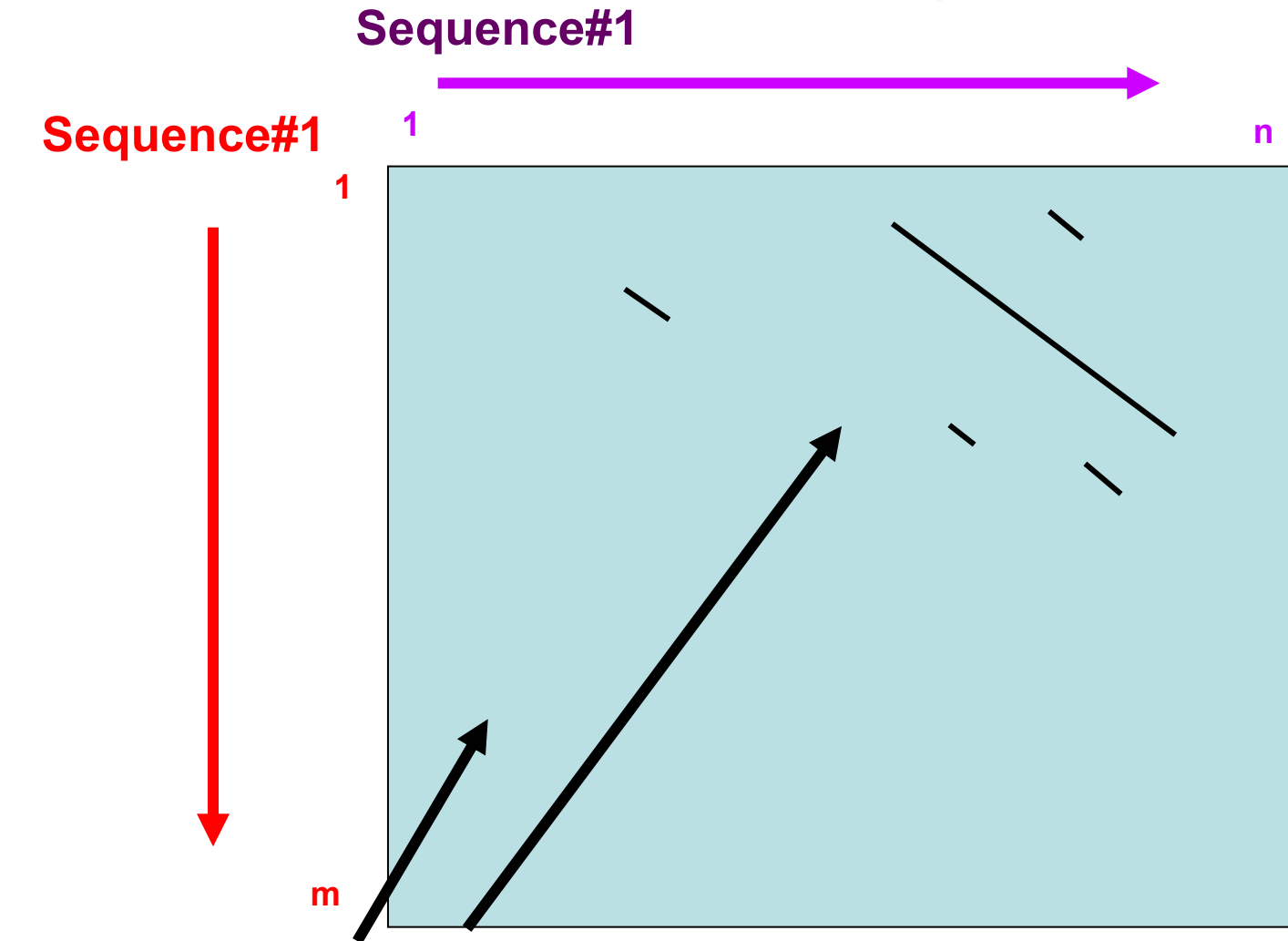
symmetric



# Now align two different sequences

- \* **Consider other similarity matrices besides identity....**
  - ***Chemical similarity*** – binary decision
  - ***Amino acid conservation*** in aligned protein families – min. similarity score (+/- window)
  - ***Average*** of multiple scoring systems

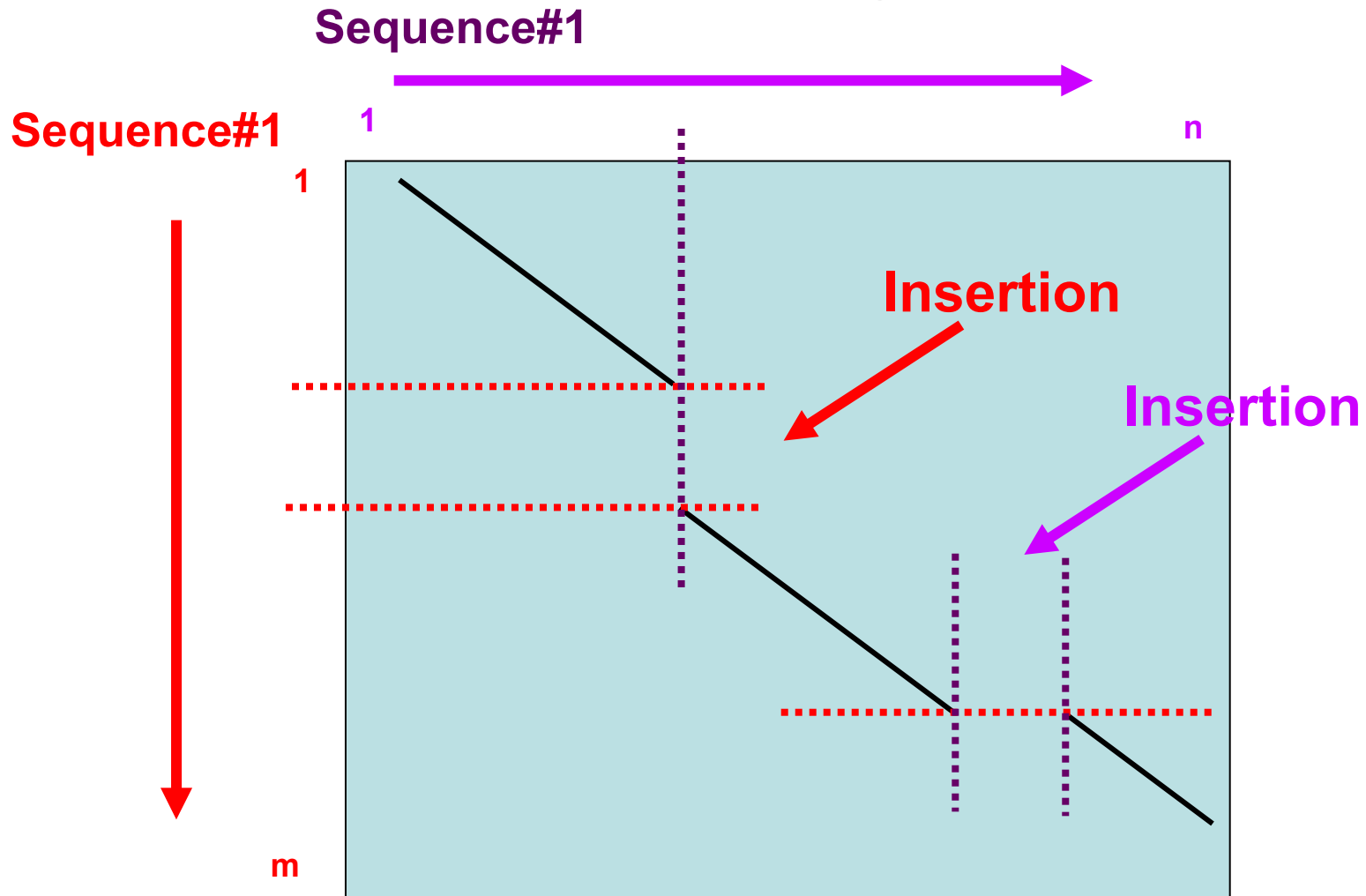
# Dot Matrix Alignments



*Note – NOT symmetric*

**A Local Alignment**

# Dot Matrix Alignments



**A Global Alignment**

# General Rules for Dot Matrices

- Advantage – let's your eyes/brain do the work – VERY EFFICIENT!!!!
- DNA comparisons – long windows and high stringencies (11/7, 15/11)
- Proteins – short windows and stringencies (1/1) EXCEPT in looking for domains – then longer windows and smaller stringencies (15/5).
- *Note – we can use these types of self-aligned matrices for more than just sequence comparisons...i.e. distance between side  $C\alpha$  atoms in 3d-structures etc....maybe later in the course!*

# Computational Efficiency

Measure efficiency in cpu run time and memory

$O()$  = “big-oh” notation

Both scale as size of the problem, measured in number of units,  $n$ , in the problem, i.e. run time is  $f(n)$ .

Analyse the asymptotic worst-case running time....

Sometimes just do the experiment and measure it....

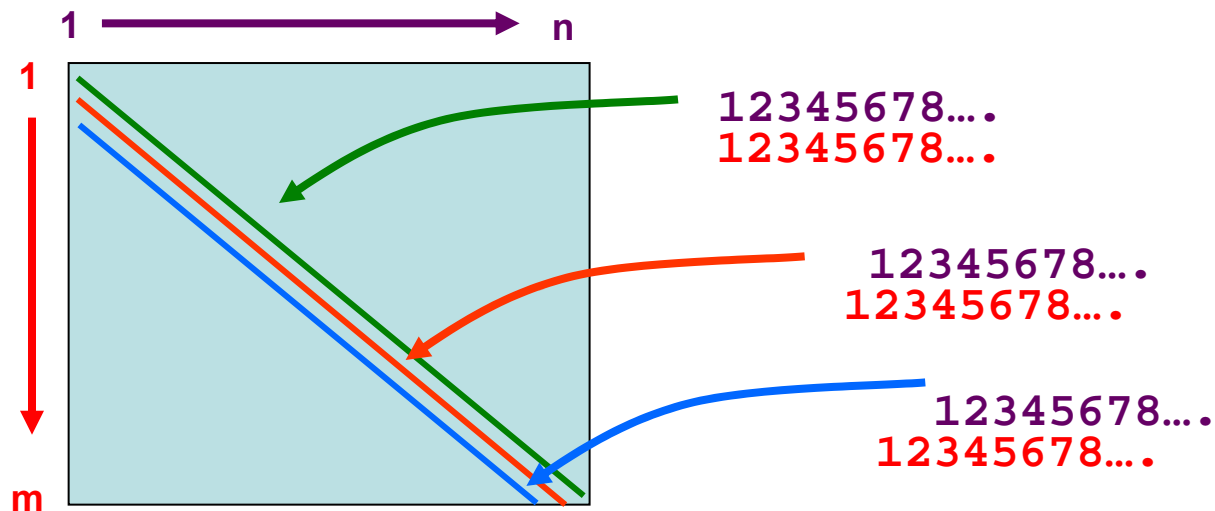
If problem scales as square of the number of units in the problem

$O(n^2)$  (=order n-squared)

## Examples

$O(n^k)$  is “polynomial time” as long as  
 $k \leq 3$  .....tractable

Consider our un-gapped dot matrix  
Global alignment:



.....essentially an  $O(mn)$  problem

## O.K. Examples

**$O(n)$  better than  $O(n \log(n))$ , better than  $O(n^2)$ , better than  $O(n^3)$**

## Terrible Examples

**$O(k^n)$  = exponential time....horrible!!!!**

**NP problems- no known polynomial time  
Solutions = non-deterministic polynomial  
Problems.**